

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

Graphs are mathematical structures that appear in many different applications of machine learning, and especially in the context of semi-supervised learning. In semi-supervised learning, we hope to exploit not only the information in the labelled examples, but also the information in the unlabelled examples by incorporating some notion of similarity between labelled and unlabelled examples. Even if our examples do not lie in a metric space, we can naturally define a notion of similarity between the examples by considering them as the vertices of a graph. One application of semi-supervised learning over graphs is in marketing and ad-tech, where we may model the flow of web traffic as a Markov Random Field. In this case, certain sites are known to be of interest to a consumer, while other sites may be of interest based on their proximity to sites previously visited on the web. Another application might be to identify suspicious accounts on a transactional network such as the block chain. In this paper we detail several different graphical models, including basic diffusion models, label propagation, regression, semi-supervised learning, and finally online learning. First, we introduce some definitions.

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 $\mathbf{A}(G)$, is the $n \times n$ matrix defined as follows. The rows and the columns of $\mathbf{A}(G)$ are indexed by $V(G)$. If $i \neq j$ then the (i, j) entry of $\mathbf{A}(G)$ 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 $\mathbf{A}(G)$ is 0 for $i = j = 1, \dots, n$. The **degree matrix** $\mathbf{D}(G)$ for $G(V, E)$ is a $n \times n$ diagonal matrix defined as

$$\mathbf{D}(G)_{i,j} := \begin{cases} d(v_i) & \text{if } i = j \\ 0 & \text{otherwise.} \end{cases} \quad (1.1)$$

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** $\mathbf{W}(G)$ of the graph $G(V, E, W)$ is defined by

$$\mathbf{W}_{ij} := \begin{cases} w(e) & \text{if } e = (i, j) \in E \\ 0 & \text{otherwise.} \end{cases} \quad (1.2)$$

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

$$\mathbf{W}(G)_{ij} = e^{-\frac{\|v_i - v_j\|^2}{2\sigma^2}}. \quad (1.3)$$

*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

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

$$D(G)_{i,i} := \sum_j W(G)_{ij} \quad (1.4)$$

1.2 Graph Laplacian and Heat Equation

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

$$L(G) := I - D^{-\frac{1}{2}} W D^{-\frac{1}{2}}, \quad (1.5)$$

and the **unnormalized graph Laplacian** is

$$L(G) := D - W. \quad (1.6)$$

Let us consider an example to understand the graph Laplacian of the graph that is shown in figure 1.1. Suppose $f : V \rightarrow \mathbb{R}$ is a real-valued function on the set of the vertices of graph $G(V, E)$ such that it

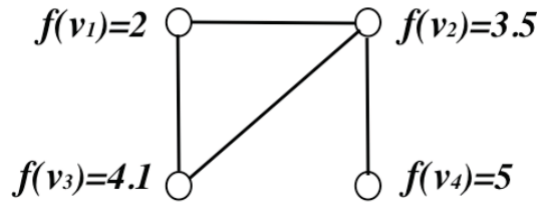


Figure 1.1: Real-Valued Functions on a Graph

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

$$A = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}. \quad (1.7)$$

Hence, the eigenvectors of the adjacency matrix, $Ax = \lambda x$, can be viewed as eigenfunctions $Af = \lambda f$. The adjacency matrix can be viewed as an operator

$$g = Af \\ g(i) = \sum_{i \sim j} f(j), \quad (1.8)$$

and it can also be viewed as a quadratic form

$$f^T A f = \sum_{e_{ij}} f(i) f(j). \quad (1.9)$$

Assume that each edge in the graph have an arbitrary but fixed orientation, which is shown in figure 1.2. Then the incidence matrix of the graph is

$$Q = \begin{bmatrix} -1 & 1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix}. \quad (1.10)$$

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

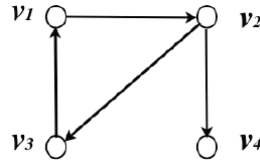


Figure 1.2: Orientation of the Graph

$$\begin{bmatrix} f(2) - f(1) \\ f(1) - f(3) \\ f(3) - f(2) \\ f(4) - f(2) \end{bmatrix} = \begin{bmatrix} -1 & 1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix} \begin{bmatrix} f(1) \\ f(2) \\ f(3) \\ f(4) \end{bmatrix}. \quad (1.11)$$

If we let

$$\mathbf{L} = \mathbf{Q}^T \mathbf{Q}, \quad (1.12)$$

then we have

$$(\mathbf{L}\mathbf{f})(v_i) = \sum_{v_i \sim v_j} [f(v_i) - f(v_j)]. \quad (1.13)$$

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

$$\mathbf{L} = \mathbf{D} - \mathbf{A} = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 3 & -1 & -1 \\ -1 & -1 & 2 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix}, \quad (1.14)$$

where the degree matrix \mathbf{D} is

$$\mathbf{D} = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}. \quad (1.15)$$

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

$$(\mathbf{L}\mathbf{f})(v_i) = \sum_{v_i \sim v_j} w_{ij} [f(v_i) - f(v_j)]. \quad (1.16)$$

Its quadratic form is

$$\mathbf{f}^T \mathbf{L} \mathbf{f} = \frac{1}{2} \sum_{e_{ij}} w_{ij} [f(v_i) - f(v_j)]^2. \quad (1.17)$$

The intuition behind a Laplacian matrix is the following. If, for instance, we apply the Laplacian operator of formula 1.14 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

$$(\mathbf{L}\mathbf{f})(v_i) = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 3 & -1 & -1 \\ -1 & -1 & 2 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix} \begin{bmatrix} f(v_1) \\ f(v_2) \\ f(v_3) \\ f(v_4) \end{bmatrix}. \quad (1.18)$$



For simplicity, let us only look at the first element

$$\begin{aligned}
 (\mathbf{L}\mathbf{f})(v_i)_1 &= \begin{bmatrix} 2 & -1 & -1 & 0 \end{bmatrix} \begin{bmatrix} f(v_1) \\ f(v_2) \\ f(v_3) \\ f(v_4) \end{bmatrix} \\
 &= 2f(v_1) - f(v_2) - f(v_3) \\
 &= -[f(v_2) - 2f(v_1) + f(v_3)] \\
 &= -[f(v_2) - f(v_1) - f(v_1) + f(v_3)]
 \end{aligned} \tag{1.19}$$

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

$$(\mathbf{L}\mathbf{f})(v_i)_1 = -[f_{k+1} - 2f_k + f_{k-1}]. \tag{1.20}$$

We recall that the second order derivative can be approximated by

$$\begin{aligned}
 f''(x) &= \frac{\frac{f(x+\Delta x) - f(x)}{\Delta x} - \frac{f(x) - f(x-\Delta x)}{\Delta x}}{\Delta x} \\
 &= \frac{f(x+\Delta x) - 2f(x) + f(x-\Delta x)}{(\Delta x)^2}.
 \end{aligned} \tag{1.21}$$

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.

We can get more intuition about the basic objects in Graph learning by using the concept of conductance. As we saw in the previous section, the Laplacian is closely related to a second derivative, and it is no coincidence that the continuous version of the Laplacian appears in the heat equation. Let \mathbf{f} be some real-valued function defined over the vertices of the graph. If we think of \mathbf{f} as a hat distribution, then we can think of $\mathbf{L}\mathbf{f}$ roughly as the flux induced at each of the vertices by that distribution over the graph, based on the graph structure. 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, then $\langle \mathbf{f}, \mathbf{g} \rangle := \mathbf{f}^T \mathbf{L}\mathbf{g}$ 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 specialized to the graph setting. In particular, we can see that the quadratic form $\mathbf{f}^T \mathbf{L}\mathbf{f}$ can be thought of as a standard metric for the “total flux” induced over the graph by some initial heat distribution \mathbf{f} . In other words, $\mathbf{f}^T \mathbf{L}\mathbf{f}$ is a measure of the smoothness of the function \mathbf{f} relative to the structure of the graph defined by \mathbf{L} .

1.3 Label Propagation

Given the graph $G(V, E, W)$, and an initial state of some real-valued function $\mathbf{f} : V \rightarrow \mathbb{R}$ defined over the vertices of the graph, we may be interested in propagating this function over the graph at time steps $t = 0, 1, 2, \dots, T$ that is consistent with the graph structure. In the simplest possible case, we might be provided with the initial state \mathbf{f}_0 and wish to see how these values propagate over the graph, as in a heat diffusion. While very simple, even scenario could be useful in settings such as fraud detection over a graph of accounts that are connected by their transactions. For example, we may have received knowledge that a small number of vertices in the graph are fraudulent, and we want to perform a quick analysis to see to what degree other vertices in the graph might be exposed to the risk emanating from these fraudulent accounts.

We proceed by deriving an update equation starting with two basic assumptions: First, we want our

function \mathbf{f}_{t+1} at each time step to be “close” to \mathbf{f}_t . Second, we want to impose the condition that \mathbf{f}_{t+1} be a “smooth” function over the graph, in the sense of minimal total flux described above. Thus we seek to minimize

$$C(\mathbf{f}_{t+1} | \mathbf{f}_t) = (\mathbf{f}_{t+1} - \mathbf{f}_t)^T (\mathbf{f}_{t+1} - \mathbf{f}_t) + \alpha \mathbf{f}_{t+1}^T \mathbf{L} \mathbf{f}_{t+1}, \quad (1.22)$$

where \mathbf{L} is the non-normalized Laplacian, and α is a constant measuring the conductivity of the graph. Taking the derivative and setting it to zero yields

$$\mathbf{f}_{t+1} = \mathbf{f}_t - \alpha \mathbf{L} \mathbf{f}_t. \quad (1.23)$$

As an example, if we were to begin with the vector of initial values $\hat{\mathbf{Y}}^0 = \mathbf{e}_1$ representing a unit point density, the above recurrence relation would show how that density propagates across the graph. We could alternatively consider the update equation to be a discretization of a system evolving in continuous-time. Note that the above equation can be viewed as a discrete difference equation

$$\mathbf{f}_{t+1} - \mathbf{f}_t = \delta \mathbf{f}_t = -\alpha \mathbf{L} \mathbf{f}_t. \quad (1.24)$$

Taking the continuous limit of this diffusion yields

$$\frac{d\mathbf{f}_t}{dt} = -\alpha \mathbf{L} \mathbf{f}_t, \quad \mathbf{f}_t(0) = \mathbf{f}_0. \quad (1.25)$$

The solution to this ODE is given by

$$\mathbf{f}_t = e^{-\alpha \mathbf{L} t} \mathbf{f}_0. \quad (1.26)$$

We can now use this continuous system in our discretization in order to derive a discrete update formula

$$\mathbf{f}_t = e^{-\alpha \mathbf{L} t} \mathbf{f}_0. \quad (1.27)$$

This system differs from the one described by the original difference equation in several ways. First, there is a closed-form solution for the value of each node at any point in time, rather than a recurrence relation, which should make the time-complexity much more efficient. But it turns out that what is gained in time efficiency is more than lost in much higher space complexity, and thus is not practical in applications with very large graphs. Consider the case where our graph contains 100000 vertices, and all vertices are connected by a single path, with the degree of each node exactly equal to 2. In this case the regular Laplacian matrix is extremely sparse, as only 0.00001 of the entries are non-zero. Using data structures designed specifically for extremely sparse matrices, storing this matrix and performing basic computations would be quite easy. On the other hand, the matrix $e^{-\alpha \mathbf{L} t}$ has strictly positive entries for all t , and therefore would require 100000 times as much memory to store.

The underlying reason for this discrepancy is that whereas the continuous Laplacian establishes a dependency among all vertices within a given connected component of the graph, the discrete Laplacian $\mathbf{L}(G)$ only propagates to immediate neighbors, as any vertices v_i and v_j , which are not connected by an edge will have value zero. To gain more intuition about this discrepancy, we recall that the matrix exponential is defined by its Taylor series

$$e^{-\alpha \mathbf{L} t} = \sum_{i=0}^{\infty} \frac{(-\alpha \mathbf{L} t)^i}{i!}. \quad (1.28)$$

Replacing our initial condition with \mathbf{f}_t and final state with \mathbf{f}_{t+1} and taking the first order approximation to this series gives

$$e^{-\alpha \mathbf{L} t} \approx \mathbf{1} - \mathbf{L}. \quad (1.29)$$

Plugging this approximation into our update formula for the continuous case yields

$$\mathbf{f}_{t+1} = e^{-\alpha \mathbf{L}} \mathbf{f}_t \approx (\mathbf{1} - \mathbf{L}) \mathbf{f}_t = \mathbf{f}_t - \alpha \mathbf{L} \mathbf{f}_t, \quad (1.30)$$



which matches exactly our formula in the discrete case. If we were to take the second order approximation, each step in our iteration would take into account not only immediate neighbours, but also neighbours of neighbours and so on as more terms from the Taylor series are kept. This in turn would reduce the total number of iteration needed to simulate the evolution of the system, so that we have achieved some of the time-efficiency of the continuous model. However, we need to pay for this more favourable time-complexity by computing and storing the matrix L^2 . Even if L is sparse, L^2 will have exponentially more non-zero entries, and will thus be much larger, and will also take much longer to compute. Therefore, if we want to optimize label propagation in practice, we may want to take as many term of the Taylor expansion as we can without exceeded our memory.

1.4 Sustained Effects of Initial State

Now suppose that we wish to maintain the effects of our initial function over time. Then we may add a term to the objective function of the form $\mathbf{f}_{t+1}^T \mathbf{f}_0$. This inner product measures the similarity between the updated function \mathbf{f}_{t+1} and the initial state vector \mathbf{f}_0 . For this problem, we will also replace the non-normalized Laplacian L with the normalized Laplacian \hat{L} , which will become useful when deriving the steady state solution below. Thus our objective function now has the form

$$C(\mathbf{f}_{t+1} | \mathbf{f}_t) = (\mathbf{f}_{t+1} - \mathbf{f}_t)^T (\mathbf{f}_{t+1} - \mathbf{f}_t) + \mathbf{f}_{t+1}^T (\mathbf{I} - \alpha \hat{L}) \mathbf{f}_{t+1} + \beta \mathbf{f}_{t+1}^T \mathbf{f}_0. \quad (1.31)$$

Straightforward calculus and linear algebra give the solution

$$\mathbf{f}_{t+1} = \alpha \hat{L} \mathbf{f}_t + \beta \mathbf{f}_0. \quad (1.32)$$

The update equation derived hence agrees with that of chapter II of Bengio, Delalleau and Roux (2006). Following Bengio, Delalleau and Roux (2006), this recurrence relation can equivalently write the entire recurrence relation in terms of $\hat{\mathbf{Y}}^0$

$$\mathbf{f}_{t+1} = (\alpha \hat{L})^t \mathbf{f}_0 + \beta \sum_{i=0}^t (\alpha \hat{L})^i \mathbf{f}_0. \quad (1.33)$$

Please note that in order to write this equation in absolute form, reducing the number of computations, we need to pay for this by having to store powers of the matrix L . It can be easily shown that the normalized Laplacian has an operator norm less than unity, $\|\hat{L}\| < 1$. Therefore, the power series $\sum_{i=0}^t (\alpha \hat{L})^i$ converges, and is equal to $(\mathbf{I} - \alpha \hat{L})^{-1}$. In addition, this implies that $\lim_{t \rightarrow \infty} (\alpha \hat{L})^t = 0$. Thus

$$\mathbf{f}_\infty = \beta (\mathbf{I} - \alpha \hat{L})^{-1} \cdot \mathbf{f}_0. \quad (1.34)$$

We will see in the sequel that the inverse or generalized inverse of the Laplacian and closely related matrices play a key role for online learning over graphs. The above is a linear system of equations which can be solved iteratively, as explained in section 2.

1.5 Regression on Graphs

Now that we have introduced the basic building blocks for graphical learning, we can begin to analyse graphical learning algorithms that are used in practice. Like in the previous sections, we begin by defining some cost function for our function f , and then we proceed to derive an update equation that minimizes that cost function. The sequel deals with algorithms that are actually used by machine learning practitioners. In this setting, we are given some set of labels some set of labels Y defined over the V , and we are asked to make predictions of those labels, $\hat{\mathbf{Y}}$.

Given a labelling \mathbf{Y} , we wish to find a set of predicted labels $\hat{\mathbf{Y}}$ that closely approximate the true labels, with some additional constraints that are familiar from the previous section. consistency with the initial

labelling can be measured by

$$\sum_{i=1}^l (\hat{y}_i - y_i)^2 = \|\hat{\mathbf{Y}}_l - \mathbf{Y}_l\|^2. \quad (1.35)$$

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

$$\frac{1}{2} \sum_{i,j=1}^n \mathbf{W}_{ij} (\hat{y}_i - \hat{y}_j)^2. \quad (1.36)$$

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 formula 1.35 and formula 1.36, which is shown in algorithm 1.1 and algorithm 1.2. The paper assumed that $G(V, E, W)$ is connected and that the vertices of the graph are numbered, and only partial information, $f(x_i) = y_i$, $1 \leq i \leq k$, 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

$$\tilde{\mathbf{y}} = (y_1 - \bar{y}, y_2 - \bar{y}, \dots, y_k - \bar{y}), \quad (1.37)$$

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

In the algorithm 1.1, \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

$$\begin{aligned} \mathcal{L} &= \frac{1}{k} \sum_i (f_i - \tilde{y}_i)^2 + \gamma \mathbf{f}^T \mathbf{S} \mathbf{f} - \mu \left(\sum f_i - 0 \right) \\ &= \frac{1}{k} \mathbf{f}^T \mathbf{f} - \frac{2}{k} \mathbf{f}^T \tilde{\mathbf{y}} + \frac{1}{k} \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} + \gamma \mathbf{f}^T \mathbf{S} \mathbf{f} - \mu \mathbf{f}^T \mathbf{1} \\ &= \mathbf{f}^T \left(\frac{1}{k} \mathbf{I} + \gamma \mathbf{S} \right) \mathbf{f} - \frac{2}{k} \mathbf{f}^T (\tilde{\mathbf{y}} + \mu \mathbf{1}) + \frac{1}{k} \tilde{\mathbf{y}}^T \tilde{\mathbf{y}}. \end{aligned} \quad (1.38)$$

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

$$\tilde{\mathbf{f}} = (\mathbf{I} + k\gamma \mathbf{S})^{-1} (\tilde{\mathbf{y}} + \mu \mathbf{1}). \quad (1.39)$$

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

$$s : \mathbf{f} \rightarrow \sum_i f_i. \quad (1.40)$$

Since s is linear, we obtain

$$0 = s(\tilde{\mathbf{f}}) = s((\mathbf{I} + k\gamma \mathbf{S})^{-1} \tilde{\mathbf{y}}) + s((\mathbf{I} + k\gamma \mathbf{S})^{-1} \mathbf{1}). \quad (1.41)$$

Therefore we can write

$$\mu = -\frac{s((\mathbf{I} + k\gamma \mathbf{S})^{-1} \tilde{\mathbf{y}})}{s((\mathbf{I} + k\gamma \mathbf{S})^{-1} \mathbf{1})} \quad (1.42)$$

In the algorithm 1.2, \mathbf{S} is a smoothness matrix, e.g. $S = L$ or $S = L^p$, $p \in \mathbb{N}$. However, here we are not



Algorithm 1.1 Tikhonov Regularization with Parameter $\gamma \in \mathbb{R}$

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

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

allowing multiple vertices in the sample. We partition \mathbf{S} as

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_1 & \mathbf{S}_2 \\ \mathbf{S}_2^T & \mathbf{S}_3 \end{bmatrix}, \quad (1.43)$$

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

$$\tilde{\mathbf{f}} = \mathbf{S}_3^{-1} \mathbf{S}_2^T \left((\tilde{y}_1, \dots, \tilde{y}_k)^T + \mu \mathbf{1} \right), \quad (1.44)$$

and

$$\mu = -\frac{s(\mathbf{S}_3^{-1} \mathbf{S}_2^T \tilde{\mathbf{y}})}{s(\mathbf{S}_3^{-1} \mathbf{S}_2^T \mathbf{1})} \quad (1.45)$$

Algorithm 1.2 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 $f(\mathbf{x}_i) = \tilde{y}_i, 1 \leq i \leq k$

$$\tilde{\mathbf{f}} = \arg \min_{\mathbf{f}=(\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_k, f_{k+1}, \dots, f_n)} \mathbf{f}^T \mathbf{S} \mathbf{f}$$

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

$$R_k(f) = \frac{1}{k} \sum_{i=1}^k (f(\mathbf{x}_i) - y_i)^2 \quad (1.46)$$

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

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

$$\mathbb{P}(|R_k(f_T) - R(f_T)| > \epsilon + \beta) \leq 2 \exp \left(-\frac{k\epsilon^2}{2(k\beta + K + M)^2} \right), \quad \forall \epsilon > 0. \quad (1.47)$$

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

$$|R_k(f_T) - R(f_T)| < \epsilon + \beta \quad (1.48)$$

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

$$\epsilon = \sqrt{\frac{2 \ln \frac{2}{\delta}}{k}} (k\beta + K + M), \quad (1.49)$$

then plug into inequality 1.48 and rearrange terms, then, with probability $1 - \delta$, we have

$$|R_k(f_T) - R(f_T)| \leq \beta + \sqrt{\frac{2 \ln \frac{2}{\delta}}{k}} (k\beta + K + M), \quad (1.50)$$

where

$$\beta = \left(\frac{3M\sqrt{tk}}{(k\gamma\lambda_1 - t)^2} + \frac{4M}{k\gamma\lambda_1 - t} \right). \quad (1.51)$$

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

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.

Building upon the previous ideas introduced in this section, we now describe a setting for semi-supervised learning over a weighted graph. We start with vertices $1, 2, \dots, l$ labelled l with their known label 1 or -1 and nodes $l+1, \dots, n$ labelled with 0. 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. In addition, we wish to take into consideration that the correct labels \mathbf{Y}_l are known for some subset of vertices.

The new semi-supervised learning objective - that we wish to minimize the difference between predicted labels and those that are known in the labelled set - can be reflected in terms of the cost function $\|\mathbf{S}\hat{\mathbf{Y}}_v - \mathbf{S}\mathbf{Y}_v\|_2^2$, where $\mathbf{S} = \mathbf{I}_{[l]}$ is a selection matrix that is equal to the identity for the first l rows and columns, and equal to 0 everywhere else. In addition, since we no longer have access to the true labels for every vertex, we add a regularization term, so that we can deal with degenerate situations, such as connected components of the graph that have no labelled vertices. Thus we seek to minimize the objective function Zhou et al. (2004)

$$C(\hat{\mathbf{Y}}) = \|\mathbf{S}\hat{\mathbf{Y}} - \mathbf{S}\mathbf{Y}\|_2^2 + \mu \hat{\mathbf{Y}}^T \mathbf{L}_w \hat{\mathbf{Y}} + \epsilon \mu \|\hat{\mathbf{Y}}\|_2^2. \quad (2.1)$$

We can use straightforward calculus techniques to minimize this expression

$$\frac{\partial C(\hat{\mathbf{Y}})}{\partial \hat{\mathbf{Y}}} = 2\mathbf{S}(\hat{\mathbf{Y}} - \mathbf{Y}) + 2\mu \mathbf{L} \hat{\mathbf{Y}} + 2\epsilon \mu \mathbf{I} \hat{\mathbf{Y}} = 0. \quad (2.2)$$

Bengio, Delalleau and Roux (2006) proposed a solution to this system of linear equations based on the Jacobi iterative method for linear systems. The Jacobi iteration algorithm gives a solution to the linear system

$$\mathbf{M}\mathbf{x} = \mathbf{b}. \quad (2.3)$$

By expanding the system as $\sum_{j=1}^n m_{ij}x_j = b_j$. Suppose we wish to solve for x_j while assuming that all other values of x remain fixed. We decompose the sum as

$$\sum_{j=1}^n m_{ij}x_j = m_{ij}x_j + \sum_{k \neq j}^n m_{ik}x_k = b_j \rightarrow x_j = b_j - \frac{\sum_{k \neq j}^n m_{ik}x_k}{m_{ij}}. \quad (2.4)$$



We can write this in matrix form as

$$\mathbf{x}^k = \mathbf{D}^{-1} \mathbf{R} \mathbf{x}^{k-1} + \mathbf{D}^{-1} \mathbf{b}, \quad (2.5)$$

where \mathbf{D} is the diagonal component of \mathbf{M} , and \mathbf{R} is the remainder, (i.e. $\mathbf{D} + \mathbf{R} = \mathbf{M}$). In our case, $\mathbf{M} = \mathbf{S} + \mu \mathbf{L} + \epsilon \mu \mathbf{I}$, and $\mathbf{b} = \mathbf{S} \mathbf{Y}$. The Jacobi method is particularly well suited for the graphical setting, since $\mathbf{M} = \mathbf{S} + \mu \mathbf{D} + \epsilon \mu \mathbf{I}$, and $\mathbf{R} = \mu \mathbf{W}$. A sufficient condition for the convergence of the Jacobi method is when the matrix \mathbf{M} is strictly diagonally dominant, meaning that the absolute value of each diagonal entry is greater than any other entry in the same row or column. Since the diagonal entries of the Laplacian are by definition positive and at least as great as any entry in the same row or column, adding the positive diagonal matrix $\mathbf{S} + \epsilon \mu \mathbf{I}$ to $\mu \mathbf{L}$ makes $\mathbf{M} = \mathbf{S} + \mu \mathbf{L} + \epsilon \mu \mathbf{I}$ strictly diagonally dominant, as desired.

Algorithm 2.3 Jacobi Iterative Label Propagation Algorithm

Compute weight matrix \mathbf{W} from formula 1.3 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 $\hat{y}_i^{(\infty)}$

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

As in the previous examples, we consider the set of real-valued functions over the vertices of the graph

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

We define the function $\langle \cdot; \cdot \rangle$ as

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

\mathbf{L} is symmetric and positive semi-definite, which suggests that by appropriately constraining our vector space containing \mathbf{f} , we could ensure the strict positivity of $\langle \mathbf{f}, \mathbf{g} \rangle$, which would imply that that our space is in fact a Hilbert space. 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 on each connected component. We can see this by noting that

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

This is because

$$\begin{aligned}
 \|g\|^2 &= \langle g, g \rangle \\
 &= \sum_{i,j=1}^n g_i L_{ij} g_j \\
 &= \sum_{i=j}^n g_i L_{ij} g_j + \sum_{i \neq j}^n g_i L_{ij} g_j \\
 &= \sum_{i=1}^n g_i^2 D_{ii} + \sum_{i \neq j}^n g_i A_{ij} g_j \\
 &= \sum_{i,j=1}^n g_i^2 A_{ij} + \sum_{i \neq j}^n g_i L_{ij} g_j \\
 &= \sum_{i,j \in E(G)} g_i^2 - 2g_i g_j \\
 &= \sum_{(i,j) \in E(G)} (g_i - g_j)^2.
 \end{aligned} \tag{3.4}$$

Note the diagonal component of L is D , the off-diagonal component of L is $-A$, $D_{ii} = \sum_{j=1}^n A_{ij}$, and in the last line of the proof we drop 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 $g_i = g_j$ for all vertices connected by an edge, i.e. when 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, $\langle f, g \rangle$ is indeed an inner product, and therefore induces a true norm $\|g\|$.

Now we can define the objective function

$$\min_{f \in H} \{ \|f\| : f_i = y_i, i = 1, \dots, l \}. \tag{3.5}$$

Note that this is very similar to the objective function as that of section 2, but our norm is being induced by a different inner product. Once L^+ has been computed, this formulation will allow us to solve a system of l equations rather than a system of n equations, as in the formulation of section 2. Since H is a Hilbert space, we can use the Riesz Representation theorem to conclude that the evaluation functional of f at any vertex, being a linear functional, can be represented as an inner product

$$f_i = f(v_i) = \langle K_i, f \rangle = K_i L f. \tag{3.6}$$

The pseudo-inverse L^+ satisfies this property for the reproducing kernel K , which can be verified by noting that

$$e_i^T L_i^+ L f = L_i^+ L f_i = f_i. \tag{3.7}$$

We call L^+ the reproducing kernel of the Hilbert space H .

Now, we return to the optimization problem 3.5. After discussing with the first author Mark Herbster, I found that in fact the formula given in that paper is slightly off. He claims that we can use the reproducing kernel property to express any f_i as $f_i = \sum_{j=1}^l K_{ij} c_j$ for some constants c . In fact, this is only true for $1 \leq i \leq l$. Since Herbster, Pontil and Wainer (2005) does not give a proof for this in the paper (or any of his other papers, or for almost any of the other claims in the paper for that matter) we give our own derivation of the minimizer of 3.5, following the arguments of Aronszajn (1950) and Wahba (1990).

Recall that the reproducing kernel of the direct sum of RKHS is equal to the sum of the reproducing kernels of the two spaces Aronszajn (1950). In our case, we consider the two subspaces (both Hilbert spaces) of H , H_1 and H_2 , where H_1 is the span of the first l basis vectors of H , corresponding to the labelled vertices, while H_2 is the span of the last $n - l$ basis vectors of H , so that $H = H_1 \oplus H_2$. Then we can write and $l \in H$ as a unique decomposition $l = l_1 + l_2$ with $l_1 \in H_1, l_2 \in H_2$. By the reproducing kernel



property, we can express any \mathbf{f}_i as $\mathbf{f}_i = \langle \mathbf{K}_i, \mathbf{f} \rangle = \langle \mathbf{K}_i, \mathbf{f}_1 \rangle + \langle \mathbf{K}_i, \mathbf{f}_2 \rangle$. Since the last $n - l$ terms of \mathbf{f}_1 and the first l terms of \mathbf{f}_2 are 0, we can equivalently write this as $\mathbf{f}_i = \langle \mathbf{K}_i, \mathbf{f} \rangle = \langle \mathbf{K}_i^1, \mathbf{f}_1 \rangle + \langle \mathbf{K}_i^2, \mathbf{f}_2 \rangle$, where $\mathbf{K}_i^1 = (\mathbf{K}_i)_{j=1}^l$ and $\mathbf{K}_i^2 = (\mathbf{K}_i)_{j=l+1}^n$. Thus \mathbf{K}^1 and \mathbf{K}^2 are the reproducing kernels of H_1 and H_2 respectively.

Since the columns of \mathbf{K}^1 and \mathbf{K}^2 form a basis for H_1 and H_2 respectively, it follows that we can express any \mathbf{f} as $\mathbf{f} = \sum_{j=1}^l \mathbf{c}_j \mathbf{K}_j^1 + \sum_{j=l+1}^n \mathbf{d}_j \mathbf{K}_j^2$. Note that the vectors \mathbf{K}_j^2 in the second sum are all perpendicular to the \mathbf{K}_j^1 , and therefore for all $1 \leq i \leq l$, we have $\mathbf{f}_i = \langle \mathbf{K}_i, \mathbf{f} \rangle = \left\langle \mathbf{K}_i, \sum_{j=1}^l \mathbf{c}_j \mathbf{K}_j^1 + \sum_{j=l+1}^n \mathbf{d}_j \mathbf{K}_j^2 \right\rangle = \left\langle \mathbf{K}_i, \sum_{j=1}^l \mathbf{c}_j \mathbf{K}_j^1 \right\rangle$.

It follows that for all $1 \leq i \leq l$, $\mathbf{f}_i = \sum_{j=1}^l \mathbf{c}_j \mathbf{K}_j^1 = \sum_{i=1}^n e_i \sum_{j=1}^l \mathbf{c}_j \mathbf{K}_{ij}^1$, or equivalently $\mathbf{f}_i = \sum_{j=1}^l \mathbf{c}_j \mathbf{K}_{ij}^1$, which matches the result of Herbster, Pontil and Wainer (2005) for $1 \leq i \leq l$. It will turn out that for the sake of minimizing 3.5, we only need to represent the first l components of \mathbf{f}_i anyway, so this is not a real oversight. In vector form, we can write have $\mathbf{f} = \mathbf{K}^1 \mathbf{c}$, where we are now referring to \mathbf{K}^1 as an $n \times r$ matrix, since the last $n - l$ columns are all 0 anyway. We can now write the Lagrange dual of 3.5 as

$$\begin{aligned} \min_{\mathbf{f} \in H} \left\{ \|\mathbf{f}\| + \lambda \sum_{i=1}^l (y_i - \mathbf{f}_i) \right\} &= \min_{\mathbf{c} \in H_1} \left\{ \|\mathbf{K}^1 \mathbf{c}\| + \lambda \sum_{i=1}^l (y_i - \mathbf{K}_i^1 \mathbf{c}) \right\} \\ &= \min_{\mathbf{c} \in H_1} \left\{ \mathbf{c}^T (\mathbf{K}^1)^T \mathbf{L} \mathbf{K}^1 \mathbf{c} + \lambda \mathbf{1}^T (\mathbf{y} - \tilde{\mathbf{K}} \mathbf{c}) \right\}, \end{aligned} \quad (3.8)$$

where $\tilde{\mathbf{K}}$ is the upper left block of \mathbf{K}^1 of dimension $l \times l$. Taking the gradient with respect to \mathbf{c} , we have $(\mathbf{K}^1)^T \mathbf{L} \mathbf{K}^1 \mathbf{c} = \tilde{\mathbf{K}} \mathbf{c} = \lambda \mathbf{1}^T \tilde{\mathbf{K}}$. Plugging in our constraint that $\mathbf{y} = \tilde{\mathbf{K}} \mathbf{c}$ we have

$$\mathbf{y} = \lambda \mathbf{1}^T \tilde{\mathbf{K}}. \quad (3.9)$$

This implies that $\mathbf{y} = \tilde{\mathbf{K}} \mathbf{c}$ and finally that $\mathbf{c} = (\tilde{\mathbf{K}})^{-1} \mathbf{y}$, which matches the solution given by Herbster, Pontil and Wainer (2005). Finally, we can express the minimizer of the original objective as

$$\mathbf{f} = \mathbf{K}^1 \mathbf{c} = \mathbf{K}^1 (\tilde{\mathbf{K}})^{-1} \mathbf{y}.$$

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 (we know Stokes' Theorem holds on general smooth manifolds, so this just a particularly simple proof of a special case).

Let $G(V, E, W)$ be a graph imbued with the topology generated by the open neighbourhoods $N_i = \{v_j : (i, j) \in E\}$. Suppose we have some subset $U(G)$ of vertices set $V(G)$, i.e. $U(G) \subset V(G)$. Then we can define the boundary of $U(G)$ as the set of all vertices contained in the closure of $U(G)$ but not contained in $U(G)$ itself $\partial U(G) = \{v_j : v_j \in \bar{U}(G) - U(G)\}$. Now suppose our vertices $V(G)$ are ordered such that $v_1, \dots, v_k \in U(G)$, $v_{k+1}, \dots, v_l \in \partial U(G)$, and $v_{l+1}, \dots, v_n \in V(G) - \bar{U}(G)$. Now let $\mathbf{1}_U$ be a binary vector with all 1 for the first k entries, and 0 elsewhere, and let $\mathbf{1}_{\partial U}$ be the binary vector with all 1 for entries $k+1$ through l , and 0 elsewhere. Then according to the above interpretation of the norm,

$$\mathbf{F}_U := \mathbf{1}_U \mathbf{L} \mathbf{g} \quad (3.10)$$

is the total flux induced by distribution \mathbf{g} over the interior of $U(G)$, while

$$\mathbf{F}_{\partial U} := \mathbf{1}_{\partial U} \mathbf{L} \mathbf{g} \quad (3.11)$$

is the total flux induced by \mathbf{g} as measured over the boundary of $U(G)$. Our claim is that

$$\mathbf{F}_U = -\mathbf{F}_{\partial U}. \quad (3.12)$$

We now give an informal proof. Let us denote by d_i the degree of vertex i , d_i^{int} by the number of edges connecting vertex i to points in $U(G)$, and d_i^{ext} by the number of edges connecting vertex i to points not in $U(G)$, so that

$$d_i = d_i^{int} + d_i^{ext}. \quad (3.13)$$

In particular, for any $v_i \in U$, d_i^{ext} is the number of edges connecting v_i to points in $\partial U(G)$, since by definition all other vertices are not connected to v_i . Then for the left hand side of formula 3.12, note that the only contribution to the i^{th} entry is given by the first k columns of \mathbf{L} , and a straightforward calculation shows that

$$\mathbf{1}_U \mathbf{L} \mathbf{g} = \sum_{i \in U} (d_i - d_i^{int}) = \sum_{i \in U} d_i^{ext}. \quad (3.14)$$

For the right hand side of formula 3.12, we can see that the only contribution to the sum is from entries $k+1$ to l , which another straightforward calculation shows that

$$\mathbf{1}_{\partial U} \mathbf{L} \mathbf{g} = \sum_{i \in \partial U} d_i^{ext}. \quad (3.15)$$

Thus

$$\mathbf{1}_U \mathbf{L} \mathbf{g} = -\mathbf{1}_{\partial U} \mathbf{L} \mathbf{g}, \quad \forall \mathbf{g} \in H. \quad (3.16)$$

Despite the simplicity of the proof and the fact that there is a very general form of Stoke's theorem over smooth manifolds, this graphical version of Stokes' theorem has not been seen by the authors, and could potentially be leveraged for new graphical algorithms.

To gain intuition about L^+ , we turn back to the representation theorem, namely that $\mathbf{f}(v_i) = \mathbf{L}_i^+ \mathbf{L} \mathbf{f}$. In the heat analogy, this could be interpreted as saying that the heat at each vertex can be expressed in terms of, or derived from, the flux through every vertex. So in this heat analogy, \mathbf{L} sends a heat distribution \mathbf{f} over each node to a flux through each vertex. Conversely, \mathbf{L}_i^+ sends some definition of fluxes over the graph back to some heat distribution that would have induced it. This relationship is shown in figure 3.1. Intuitively,

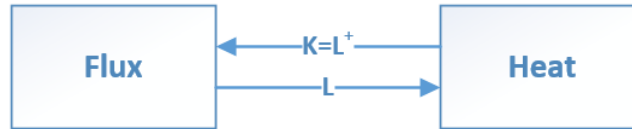


Figure 3.1: Bidirectional Mapping Between Heat and Flux

if the graph is disconnected this distribution should not be unique, since a constant added to the distribution in each of the connected components will not affect the flux induced through the graph. Mathematically, this is affirmed by noting that

$$\mathbf{L}^+ = \mathbf{L}^{-1} \quad (3.17)$$

if and only if all eigenvalues of \mathbf{L} are non-zero which, is exactly when the graph is connected.

3.2 Online Projections over Graphs

With that in mind, we can introduce the mathematical framework for the online projection algorithm of Herbster, Pontil and Wainer (2005). Herbster seeks to minimize the number of mistakes (cumulative error) on the sequence of pattern-label pairs over a series of trials. The key tool for this algorithm is projection, which is defined by definition 3.1.

Definition 3.1. The projection of a point $\mathbf{w} \in \mathcal{H}$, which is a Hilbert space with inner product $\langle \cdot, \cdot \rangle$ and norm $\|\cdot\| = \sqrt{\langle \cdot, \cdot \rangle}$, onto a closed convex non-empty set $\mathcal{N} \subseteq \mathcal{H}$ is defined by

$$P(\mathcal{N}; \mathbf{w}) := \arg \min_{\mathbf{u} \in \mathcal{N}} \|\mathbf{u} - \mathbf{w}\|. \quad (3.18)$$

From the definition, we can see that \mathbf{u} any point in \mathcal{N} that is closest to \mathbf{w} . If we assume $\mathbf{u}^* = P(\mathcal{N}; \mathbf{w}) = \arg \min_{\mathbf{u} \in \mathcal{N}} \|\mathbf{u} - \mathbf{w}\|$, then the solution is characterized by an orthogonality condition: \mathbf{u}^* is the solution if and only if

$$\langle \mathbf{w} - \mathbf{u}^*, \mathbf{u} - \mathbf{u}^* \rangle = 0. \quad (3.19)$$

Therefore by the Pythagorean Theorem, we have

$$\|\mathbf{u} - \mathbf{w}\|^2 = \|\mathbf{u} - \mathbf{u}^*\|^2 + \|\mathbf{u}^* - \mathbf{w}\|^2 \leq \|\mathbf{u} - \mathbf{w}\|^2. \quad (3.20)$$

Hence we have the following theorem.

Theorem 3.1. If \mathcal{N} is a closed convex subset of \mathcal{H} , then, for every $\mathbf{u} \in \mathcal{N}$ and $\mathbf{w} \in \mathcal{H}$, we have that

$$\|\mathbf{u} - \mathbf{w}\|^2 \geq \|\mathbf{u} - P(\mathcal{N}; \mathbf{w})\|^2 + \|P(\mathcal{N}; \mathbf{w}) - \mathbf{w}\|^2. \quad (3.21)$$

In particular, if \mathcal{N} is an affine set the equality holds.

Based on theorem 3.1, and we let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^l \subseteq \mathcal{X} \times \{-1, +1\}$ be an example set, let \mathcal{D}_t denote the first t examples of \mathcal{D} . if we consider an online learning algorithm (algorithm 3.4) such that at each iteration, \mathbf{w}_t is updated by the projection 3.18, then it is convergent. Because on any trial $t = 1, 2, \dots, l$, the theorem 3.1 implies, for all $\mathbf{u} \in \mathcal{U}$, that $\|\mathbf{u} - \mathbf{w}_t\|^2 \geq \|\mathbf{u} - \mathbf{w}_{t+1}\|^2 + \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2$, which is $\|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \leq \|\mathbf{u} - \mathbf{w}_t\|^2 - \|\mathbf{u} - \mathbf{w}_{t+1}\|^2$. Therefore, by summing it from $t = 1$ to $t = l$, we have $\sum_{t=1}^l \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \leq \|\mathbf{u} - \mathbf{w}_1\|^2 - \|\mathbf{u} - \mathbf{w}_{l+1}\|^2$, where $\mathcal{U}_t \in \mathcal{H}$, and $\mathbf{u} \in \cap_{t=1}^l \mathcal{U}_t$.

Algorithm 3.4 Prototypical Projection Algorithm

Input: A sequence of closed convex sets $\{\mathcal{U}_t\}_{t=1}^l \subset \mathcal{H}$

Initialization: $\mathbf{w}_1 \in \mathcal{H}$

for $t = 1, 2, \dots, l$ **do**

$\mathbf{w}_{t+1} = P(\mathcal{U}_t; \mathbf{w}_t)$

end

Lemma 3.1. If $(x, y) \in \mathcal{H} \times \{-1, +1\}$ and $\mathbf{w} \in \mathcal{H}$, then

$$P(\{\mathbf{u} : \langle \mathbf{u}, \mathbf{x} \rangle = y\}; \mathbf{w}) = \mathbf{w} + \frac{(y - \langle \mathbf{w}, \mathbf{x} \rangle)}{\|\mathbf{x}\|^2} \mathbf{x} \quad (3.22)$$

$$P(\{\mathbf{u} : y \langle \mathbf{u}, \mathbf{x} \rangle \geq 1\}; \mathbf{w}) = \mathbf{w} + \frac{y \max(0, 1 - y \langle \mathbf{w}, \mathbf{x} \rangle)}{\|\mathbf{x}\|^2} \mathbf{x} \quad (3.23)$$

Proof. We give a proof of formula 3.22, which is almost exactly the same as the proof of formula 3.23. We know that $\{\mathbf{u} : \langle \mathbf{u}, \mathbf{x} \rangle = y\} = \left\{ \mathbf{u} : \left\langle \mathbf{u} - \frac{y \mathbf{x}}{\|\mathbf{x}\|^2}, \mathbf{x} \right\rangle = 0 \right\}$. Now let $\mathbf{v}(\mathbf{u}) = \mathbf{u} - \frac{y \mathbf{x}}{\|\mathbf{x}\|^2}$. Then we can find this projection of \mathbf{w} onto $\{\mathbf{u} : \langle \mathbf{u}, \mathbf{x} \rangle = y\}$ by first projecting $\mathbf{v}(\mathbf{w})$ onto $\{\mathbf{u} : \langle \mathbf{u}, \mathbf{x} \rangle = 0\}$, $\mathbf{v}^* =$

$\text{Proj}(\{\mathbf{u} : \langle \mathbf{u}, \mathbf{x} \rangle = 0\}; \mathbf{v}(\mathbf{w}))$ and then taking $\mathbf{w}^* = \mathbf{v}^{-1}(\mathbf{v}^*)$.

$$\begin{aligned}
\mathbf{v}^* &= \text{Proj}(\{\mathbf{u} : \langle \mathbf{u}, \mathbf{x} \rangle = 0\}; \mathbf{v}(\mathbf{w})) \\
&= \mathbf{v}(\mathbf{w}) - \text{Proj}_{\mathbf{x}}(\mathbf{v}(\mathbf{w})) \\
&= \mathbf{v}(\mathbf{w}) - \frac{\langle \mathbf{v}(\mathbf{w}), \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} \mathbf{x} \\
&= \mathbf{v}(\mathbf{w}) + \frac{\left\langle \mathbf{w} - \frac{\mathbf{x}y}{\|\mathbf{x}\|^2}, \mathbf{x} \right\rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} \mathbf{x} \\
&= \mathbf{v}(\mathbf{w}) + \frac{y - \langle \mathbf{w}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} \mathbf{x}.
\end{aligned} \tag{3.24}$$

Thus

$$\mathbf{w}^* = \mathbf{w} + \frac{(y - \langle \mathbf{w}, \mathbf{x} \rangle)}{\|\mathbf{x}\|^2} \mathbf{x} \tag{3.25}$$

□

Based on the prototypical projection algorithm 3.4 and lemma 3.1, the main algorithm of Herbster, Pontil and Wainer (2005) is given in algorithm 3.5. In the algorithm, M is the set of trials in which the algorithm predicts incorrectly, U_t is the index set, which is determined by the **strategy** function on trial t . This in turn determines the **noncyclic** feasible set \mathcal{U}_t or a **cyclic** feasible sets $\{\mathcal{U}_{t_1}, \mathcal{U}_{t_2}, \dots\}$, which each feasible set \mathcal{U}_t consists of those hypothesis vectors which are “compatible” with a subset of the past examples (x_i, y_i) . In this algorithm, there are two strategies for the non-cyclic case and cyclic case. The non-cyclic strategy projects the current hypothesis \mathbf{w}_t to a single affine set $\text{af}(U) := \cap_{i \in U} \{\mathbf{u} \in \mathcal{H} : \langle \mathbf{u}, \mathbf{x}_i \rangle = y_i\}$ that is determined by the example indices in U_t . The cyclic strategy projects on a sequence of half-spaces $\text{hs}(U) := \cap_{i \in U} \{\mathbf{u} \in \mathcal{H} : y_i \langle \mathbf{u}, \mathbf{x}_i \rangle \geq 1\}$ determined by the example indices in U_t . The boolean **aggressive** controls whether an update is made on every trial (aggressive) or only on those trials when a mistake occurs (non-aggressive).

Algorithm 3.5 Online Projection Algorithm

Input: $\mathcal{D}_\ell = (x_1, y_1), \dots, (x_\ell, y_\ell) \in \mathcal{H} \times \{-1, +1\}$.

Initialization: $\mathbf{w}_1 \in \mathcal{H}$; **cyclic** = $\langle \text{TRUE} | \text{FALSE} \rangle$;

aggressive = $\langle \text{TRUE} | \text{FALSE} \rangle$

for $t = 1, \dots, \ell$ **do**

Predict: receive \mathbf{x}_t

$\hat{y}_t = \text{sign}(\langle \mathbf{w}_t, \mathbf{x}_t \rangle)$

Update: receive y_t

if $\hat{y}_t \neq y_t$ **then** $M = M \cup \{t\}$

$U_t = \text{strategy}(M, \mathbf{w}_t, \mathcal{D}_t)$

if **cyclic** **then**

$\mathbf{w}_{t_1} = \mathbf{w}_t$; $i = 1$

1 **if** $\langle \mathbf{w}_{t_i}, \mathbf{x}_t \rangle y_t \leq 0$ **or** **aggressive** **then**

$\mathbf{w}_{t_{i+1}} = P(\text{hs}(\{t\}); \mathbf{w}_{t_i})$; $i = i + 1$

while $\exists \tau_i \in U_t : \langle \mathbf{w}_{t_i}, \mathbf{x}_{\tau_i} \rangle y_{\tau_i} < 0$ **do**

2 $\mathbf{w}_{t_{i+1}} = P(\text{hs}(\{\tau_i\}); \mathbf{w}_{t_i})$; $i = i + 1$

end

$\mathbf{w}_{t+1} = \mathbf{w}_{t_i}$

else

if $\langle \mathbf{w}_t, \mathbf{x}_t \rangle y_t \leq 0$ **or** **aggressive** **then**

3 $\mathbf{w}_{t+1} = P(\text{af}(U_t); \mathbf{w}_t)$

end

end

Theorem 3.2. *If $\{(x_i, y_i)\}_{i=1}^l \subseteq \mathcal{H} \times \{-1, +1\}$ is a sequence of examples, $\mathbf{w}_1 \in \mathcal{H}$ a start vector and M the set of trials in which Algorithm 3.5 predicted incorrectly, then the cumulative number of mistakes $|M|$ of the*



algorithm is bounded by

$$|M| \leq \|\mathbf{u} - \mathbf{w}_1\|^2 B \quad (3.26)$$

for all $\mathbf{u} \in \text{hs}(\{1, \dots, l\})$ with cyclic updating and for all $\mathbf{u} \in \text{af}(\{1, \dots, l\})$ with non-cyclic updating, where

$$B = \left(\frac{1}{|M|} \sum_i (\|\mathbf{x}_i\|^2)^{-1} \right)^{-1}. \quad (3.27)$$

Theorem 3.2 bounds the number of mistakes as proportional to the squared norm (“complexity”) of the predictor \mathbf{u} . We can sketch the proof. For instance, let us consider the non-cyclic case. On each trial t there is a projection to a single affine set $\mathcal{U}_t = \text{af}(U_t)$. Since we showed that for all $\mathbf{u} \in \text{af}(\{1, \dots, l\})$, we have $\sum_{t=1}^l \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \leq \|\mathbf{u} - \mathbf{w}_t\|^2 - \|\mathbf{u} - \mathbf{w}_{l+1}\|^2$. If a mistake has occurred at trial t , $y_t = \langle \mathbf{w}_{t+1}, \mathbf{x}_t \rangle$, then $1 \leq |\langle \mathbf{w}_t, \mathbf{x}_t \rangle - y_t| \leq |\langle \mathbf{w}_t - \mathbf{w}_{t+1}, \mathbf{x}_t \rangle - y_t| \leq \|\mathbf{w}_t - \mathbf{w}_{t+1}\| \|\mathbf{x}_t\|$. Therefore, $1 \leq \|\mathbf{w}_t - \mathbf{w}_{t+1}\| \|\mathbf{x}_t\|$, which implies $\|\mathbf{x}_t\|^{-1} \leq \|\mathbf{w}_t - \mathbf{w}_{t+1}\|$. Hence, $\sum_{t \in M} \|\mathbf{x}_t\|^{-2} \leq \sum_{t=1}^l \|\mathbf{w}_t - \mathbf{w}_{t+1}\|^2$. Combining them, we have $\sum_{t \in M} \|\mathbf{x}_t\|^{-2} \leq \|\mathbf{u} - \mathbf{w}_t\|^2 + \|\mathbf{u} - \mathbf{w}_{l+1}\|^2$, which gives us $|M| \leq (\|\mathbf{u} - \mathbf{w}_t\|^2 - \|\mathbf{u} - \mathbf{w}_{l+1}\|^2) B$. After dropping the term $\|\mathbf{u} - \mathbf{w}_{l+1}\|^2$, we then get formula 3.26. The argument for the cyclic case follows the above argument except that the left hand side of the analogous inequality to $\sum_{t=1}^l \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \leq \|\mathbf{u} - \mathbf{w}_t\|^2 - \|\mathbf{u} - \mathbf{w}_{l+1}\|^2$ contains sub-terms due to repeatedly cycling through the past examples, which is line 2 in algorithm 3.5. These terms may all be lower-bounded by zero except the term corresponding to the first projection, which is line 1 in algorithm 3.5.

3.3 Online Learning over Graphs

As shown in section 3.1, the reproducing kernel of the Hilbert space defined over the graph is given by the pseudo-inverse of the Laplacian matrix, \mathbf{L}^+ . This suggests that we can apply the above projection algorithm, with its associated guarantees, by letting the input vector \mathbf{x}_t be replaced by the t^{th} row of \mathbf{L}^+ . We saw in the previous section that the number of cumulative mistakes is upper bounded by the product of two terms: one involving $\|\mathbf{u} - \mathbf{w}_1\|^2$, and the other involving the harmonic mean of the input patterns that were misclassified. Since the initial weight vector \mathbf{w} is arbitrary, we will consider the first term to be bounded more generally by $\|\mathbf{u}\|^2$. In the following analysis, we will derive upper bounds for each of these terms using the properties of our graph, and of the labellings given to us.

Theorem 3.3. *Let \mathbf{g}^+ be the positively labelled vertices, and let \mathbf{g}^- be the negatively labelled vertices. Let $|\mathbf{g}^+| = n^+$, $|\mathbf{g}^-| = n^-$, and assume WLOG that $n^+ \geq n^-$. We claim that $\|\mathbf{u}\|^2 \leq \partial(\mathbf{g}^+, \mathbf{g}^-)(1 + \frac{n^+}{n^-})^2$*

Proof. To show this, we consider the vector \mathbf{f} , with $f_i = 1$ if $g_i = 1$, and $f_i = -\frac{n^+}{n^-}$ if $g_i = -1$. Now suppose we order our vertices so that the first n^+ vertices are those in \mathbf{g}^+ , while the last n^+ vertices are those in \mathbf{g}^- . Now consider the Laplacian \mathbf{L} of this matrix with the vertices ordered in such a way. We are interested in the sums of the entries in the off-diagonal blocks, the sums of whose entries are each the negative of the number of intra-partition edges. We will denote this sum by S , as opposed to the sub-matrix \mathbf{S} .

$$\mathbf{L} = \begin{bmatrix} \mathbf{A} & -\mathbf{S} \\ -\mathbf{S}^T & \mathbf{B} \end{bmatrix}. \quad (3.28)$$

Note that the sum of entries in the left partition of the matrix for any given row is always the negative of the sum to the right of the partition. Using this and the fact that the sum of entries in both \mathbf{S} and \mathbf{S}^T are both equal to S , we find after some algebra that

$$\|\mathbf{f}\|^2 = \langle \mathbf{f}, \mathbf{f} \rangle = \mathbf{f}^T \mathbf{L} \mathbf{f} = \left(1 + \frac{n^+}{n^-}\right)^2 S \quad (3.29)$$

□

Finally, combined with the bound on the power mean of the input patterns in terms of the graph diameter, we achieve a bound on the number of mistakes.

As before, our weight vector \mathbf{g} represents a real-valued function on our Hilbert space. But in the graph setting, this Hilbert space and its reproducing kernel \mathbf{L}^+ and the vector \mathbf{g} admit a specific interpretation.

3.4 Drawbacks with the above Approach

Large diameter graphs

With that in mind, we can introduce the mathematical framework for the online projection algorithm of Herbster, Pontil and Wainwright (2005). Herbster seeks to minimize

- projection description (Tom)
- learning bounds (Paul)

Extension of the herbster paper - use different laplacians simultaneously by considering them each as different kernels in some set of kernels, and then applying the objective function for multiple kernels (need to flesh out this idea)

Bibliography

- Aronszajn, Nachman (1950). "Theory of Reproducing Kernels". In: *Transactions of the American Mathematical Society* vol. 68. no. 3, pp. 337–404. ISSN: 0002-9947 (cit. on p. [II](#)).
- Bapat, Ravindra B. (2014). *Graphs and Matrices*. London, UK: Springer-Verlag. 193 pp. ISBN: 9781447165682.
- 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. [7–9](#)).
- Bengio, Yoshua, Olivier Delalleau and Nicolas Le Roux (2006). *Label Propagation and Quadratic Criterion*. Ed. by Olivier Chapelle, Bernhard Schölkopf and Alexander Zien. Semi-Supervised Learning. Cambridge, Massachusetts, USA: The MIT Press. 528 pp. ISBN: 9780262033589 (cit. on pp. [6](#), [9](#)).
- 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 pp. [10–13](#), [15](#), [17](#)).
- 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](#)).
- Wahba, Grace (1990). *Spline Models for Observational Data*. CBMS-NSF Regional Conference Series in Applied Mathematics. Philadelphia, Pennsylvania, USA: SIAM. 161 pp. ISBN: 9780898712445 (cit. on p. [II](#)).
- 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. [9](#)).
- 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.

