

Semi-Supervised Learning on Riemannian Manifolds

University of Vienna



Peter Mühlbacher

October 1, 2014

Abstract

Contents

1	Mapping Graphs to the Real Line	3
1.1	Classification	3
1.2	An Implicit Definition of the Classifier	4
1.2.1	The Discrete Case	4
1.2.2	The Continuous Case	5
A	Implementations	6

Introduction

Chapter 1

Mapping Graphs to the Real Line

1.1 Classification

Given a space X and a probability distribution τ on $X \times \mathbb{R}^1$, a *regression problem* is the problem of finding a function $f : X \rightarrow \mathbb{R}$, $f(x) = E_{\tau_x}(y)$, where y has the distribution τ_x (the marginal distribution of x on \mathbb{R}).

In real-world applications we have a finite set of samples $(x_1, y_1), \dots, (x_n, y_n)$ from which we reconstruct a function that should

1. describe the given examples well
2. make the probability of making large errors for new examples small

This problem is closely related to the idea of interpolation and as we already know fulfilling the first requirement perfectly hardly ever goes hand in hand with the second one. In practice one wants to balance the complexity² of a function and its ability to describe given data.

Classification may be seen as a special case of regression where we aim to find a function $f : X \rightarrow \{c_1, \dots, c_n\}$.

¹another way to put this would be that for every $x \in X$ there exists a probability distribution τ_x specifying how likely this x goes together with some real number r

²Depending on the concrete implementation there are different definitions of complexity. In the theory of splines (or kernel methods in machine learning) $\|f\|_{\mathbb{H}}$, where \mathbb{H} is a reproducing kernel Hilbert space, may be seen as one; there also exist combinatorial approaches and in this work we will introduce a smoothness functional to be minimized.

1.2 An Implicit Definition of the Classifier

1.2.1 The Discrete Case

To develop a better understanding the discrete case will be dealt with in the first place.

Given a set of points $\{x_i\}_{i=1,\dots,n}$, $x_i \in \mathbb{R}^l$ a weighted graph G can be constructed. Nearby³ nodes x_i, x_j are connected with weights $W_{ij} = e^{-\frac{\|x_i - x_j\|^2}{t}}$, the weight between unconnected nodes is set to 0.

Consider the problem of mapping the weighted graph to the real line so that connected points stay together as close as possible.⁴

As the above defined weights of adjacent nodes are defined in such a way that they are indirectly proportional to their (euclidian) distance it seems natural that $(y_1, \dots, y_n)^T$ is a good map if

$$\sum_{ij} (y_i - y_j)^2 W_{ij} \quad (1.1)$$

is minimal.

Let D be the degree matrix of the graph G , then its Laplacian matrix \mathcal{L} is defined as

Definition 1.1. $\mathcal{L} := D - W$

and (1.1) may be written as

$$\begin{aligned} \sum_{ij} (y_i - y_j)^2 W_{ij} &= \sum_{ij} (y_i^2 + y_j^2 - 2y_i y_j) W_{ij} \\ &= \sum_i y_i^2 D_{ii} + \sum_j y_j^2 D_{jj} - 2 \sum_{ij} y_i y_j W_{ij} \\ &= 2y^T \mathcal{L} y \end{aligned}$$

Not only do we see that the minimization problem reduces to finding

$$\operatorname{argmin}_y y^T \mathcal{L} y$$

but also that \mathcal{L} is positive semidefinite.

However, an additional constraint ($y^T \mathbf{1} = 0$, where $\mathbf{1}$ is the constant function mapping all nodes to 1; it is easy to see that this is an eigenfunction to the eigenvalue 0) is needed as the trivial solution is not excluded.

³There are two major common approaches of determining when two nodes are adjacent: k nearest neighbours (find the k nearest nodes) and ϵ -neighbourhoods (two nodes x, y are connected if and only if $\|x - y\| < \epsilon$)

⁴For simplicity's sake assume that the graph is connected.

depending on the length of the work either remark that an explanation will be given later on (heat diffusion resembles this (=Gaussian heat kernel) or not at all (in this case set weights to 1 (which may also be seen as $t \rightarrow \infty$, all mathematical properties should stay the same))

Additionally, we want to assign some kind of importance to each of the nodes, based on how many adjacent nodes it has; we do this by adding the additional constraint $y^T D y = 1$ which, due to the nature of the objective function (1.1), leads to the generalized eigenvalue problem $\mathcal{L}y = \lambda D y$.

1.2.2 The Continuous Case

Definition 1.2. $\mathcal{L}f := -\operatorname{div} \nabla f$

explain why eigenvalues are important here (could resemble the proof why the Hölder norm for $p = 2$ induces the spectral norm(?))

define only for \mathbb{R}^n or more generally?

Appendix A

Implementations