

# Hyperspectral image asbestos classification with supervised and semi-supervised learning

Anh Van Giang (V)

June 14, 2022

## Abstract

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Dimensionality Reduction</b>	<b>2</b>
<b>3</b>	<b>Optimization Theory</b>	<b>2</b>
3.1	Problem Formulation . . . . .	2
3.2	Lagrangian Theory . . . . .	3
<b>4</b>	<b>Support Vector Machine</b>	<b>5</b>
4.1	Hard Margin SVM . . . . .	5
4.2	Soft Margin SVM . . . . .	7
<b>5</b>	<b>Positive Definite Kernels and the Reproducing Kernel Hilbert Space</b>	<b>8</b>
5.1	Positive Definite Kernels . . . . .	8
5.2	Reproducing Kernel Hilbert Space . . . . .	12
<b>6</b>	<b>Riemannian Manifold</b>	<b>14</b>
6.1	Topology prerequisites . . . . .	14
6.2	Smooth manifolds . . . . .	15
6.3	Smooth maps . . . . .	17
6.4	Tangent Vectors . . . . .	17
6.5	Submersions and Embeddings . . . . .	20
6.6	Covector fields and the differential of a function . . . . .	20
6.7	Tensors . . . . .	21
6.8	Riemannian Manifold . . . . .	22
<b>7</b>	<b>Manifold Regularization and Laplacian SVM</b>	<b>23</b>
7.1	Regularization . . . . .	24
7.2	Manifold Regularization . . . . .	24
<b>8</b>	<b>Results</b>	<b>24</b>
<b>9</b>	<b>Conclusion</b>	<b>24</b>

## Introduction

## Dimensionality Reduction

## Optimization Theory

Optimization theory is a branch of mathematics that concerns itself with characterizing the solutions of finding a set of parameters that minimizes (or maximizes) a certain cost function, usually with respect to some constraints. This section aims to provide a brief introduction of the field of optimization and the methods used to solve the Support Vector Machine (SVM) problem.

### 3.1 Problem Formulation

In the context of real-world applications not restricted to SVM, one may encounter problems that are of the form of maximizing or minimizing a function w.r.t some constraints. These problems can be formulated as:

**Definition 3.1.1.** (Primal optimization problem) Given functions  $f$ ,  $g_i$ ,  $i = 1, \dots, k$ , and  $h_i$ ,  $i = 1, \dots, m$ , defined on a domain  $\Omega \subseteq \mathbb{R}^n$ ,

$$\begin{aligned} \min \quad & f(w), \quad w \in \Omega \\ \text{s.t.} \quad & g_i(w) \leq 0, \quad i = 1, \dots, k, \\ & h_i(w) = 0, \quad i = 1, \dots, m, \end{aligned} \tag{3.1}$$

where  $f(w)$  is called the *objective* function, and the remaining relations are called, respectively, the inequality and equality constraints. The optimal value of the objective function is called the value of the optimization problem.

Notice that we can easily convert a maximization problem to a minimization one by reverting the sign of  $f(w)$ . Using the optimization problem in Definition 3.1.1, define

$$R = \{w \in \Omega \mid g_i(w) \leq 0, h_i(w) = 0\}$$

as the feasible region i.e the region of the domain where the objective function is defined and all the constraints are satisfied.

A solution of the optimization problem is a point  $w^* \in R$  such that for any  $w \in R$ ,  $f(w) > f(w^*)$ , also known as the global minimum (or maximum).

**Definition 3.1.2.** A point  $w^* \in \Omega$  is called a *local minimum* of  $f(w)$  if there exists some  $\epsilon > 0$  such that for all  $w$  that satisfies  $\|w - w^*\| < \epsilon$ ,  $f(w^*) \leq f(w)$ .

An inequality constraint  $g_i(w) \leq 0$  is said to be *active* if the solution  $w^*$  satisfies  $g_i(w^*) = 0$ , i.e, on the boundary, and it is said to be *inactive* otherwise. To transform an inequality constraint into an equality one, *slack variables* denoted as  $\xi$  can be applied as follow:

$$g_i(w) \leq 0 \iff g_i(w) + \xi_i = 0, \text{ with } \xi_i \geq 0.$$

These variables will be used extensively when we introduce Soft Margin SVM, to indicate a certain amount of "looseness" in the constraint. Most if not all the loss functions used in this paper will be convex so we will restrict the content accordingly.

**Definition 3.1.3.** A set  $\Omega \subseteq \mathbb{R}^n$  is convex, if for all  $x, y \in \Omega$ , and for all  $\theta \in [0, 1]$

$$\theta x + (1 - \theta)y \in \Omega.$$

A function  $f : \Omega \rightarrow \mathbb{R}$  is called *convex* if for all  $w, u \in \Omega$ , and for  $\theta \in [0, 1]$ ,

$$f(\theta w + (1 - \theta)u) \leq \theta f(w) + (1 - \theta)f(u).$$

If a strict inequality holds then the function is said to be *strictly convex*.

The reason why we mainly use convex functions in this paper will be clear from the following proposition.

**Proposition 3.1.1.** *If a function  $f : \Omega \rightarrow \mathbb{R}$  is convex, then any local minimum of the unconstrained optimization with the objective function  $f$  is also a global minimum.*

*Proof.* Let  $w^* \in \Omega$  be a local minimum then there exists  $\epsilon > 0$  s.t

$$f(w^*) \leq f(w), \quad \forall w \in B(w^*, \epsilon) = \{w : \|w - w^*\| < \epsilon\}.$$

Suppose that there is a  $u \in \Omega$  with

$$f(u) < f(w^*)$$

then because  $\Omega$  is convex, we have

$$\theta w^* + (1 - \theta)u \in \Omega, \quad \forall \theta \in [0, 1].$$

By the convexity of  $f$ ,

$$\begin{aligned} f(\theta w^* + (1 - \theta)u) &\leq \theta f(w^*) + (1 - \theta)f(u) \\ &< \theta f(w^*) + (1 - \theta)f(w^*) \\ &= f(w^*). \end{aligned}$$

Choose a  $\theta$  sufficiently close to 1 such that  $\theta w^* + (1 - \theta)u \in B(w^*, \epsilon)$  but then  $f(\theta w^* + (1 - \theta)u) < f(w^*)$  which is a contradiction.  $\square$

Thus, by imposing the convexity condition on our optimization problem, we can guarantee that a unique solution exists. The next part will present the Lagrange multipliers technique to solve convex quadratic optimization problem.

## 3.2 Lagrangian Theory

The methods of Lagrange multipliers and the Lagrangian function were first developed by Lagrange in 1797 [1] and extended by Kuhn, Tucker [2] in 1951 to allow for inequality constraints. These theories will provide the sufficient solutions for the Support Vector Machine problem.

**Theorem 3.2.1.** (Fermat) *A necessary condition for  $w^*$  to be a minimum of  $f(w)$ ,  $f \in C^1$ , is*

$$\frac{\partial f(w^*)}{\partial w} = 0.$$

*This condition, together with convexity of  $f$ , is also a sufficient condition.*

**Definition 3.2.1.** Given an optimization problem with objective function  $f(w)$ , and equality constraints  $h_i(w) = 0$ ,  $i = 1, \dots, m$ , we define the *Lagrangian function* as

$$L(w, \alpha) = f(w) + \sum_{i=1}^m \alpha_i h_i(w)$$

where the coefficients  $\alpha_i$  are called the *Lagrange multipliers*.

The Lagrangian function incorporates information about both the objective function and the constraints, whose stationary points can be used to find solutions.

**Theorem 3.2.2.** (Lagrange) *A necessary condition for a normal point  $w^*$  to be a minimum of  $f(w)$  subject to  $h_i(w) = 0$ ,  $i = 1, \dots, m$  with  $f, h_i \in C^1$ , is*

$$\begin{aligned} \frac{\partial L(w^*, \alpha^*)}{\partial w} &= 0, \\ \frac{\partial L(w^*, \alpha^*)}{\partial \alpha} &= 0, \end{aligned}$$

*for some values  $\alpha^*$ . The above conditions are also sufficient provided that  $L(w, \beta^*)$  is a convex function of  $w$ .*

A simple geometric interpretation of this theorem with only 1 equality constraint can be given as trying to "climb" a surface defined by  $f(w)$  subjected to a path  $h(w) = 0$ . Then, the solution would be the point where the gradient of the surface is parallel to the gradient of the curve,  $\nabla f = \lambda \nabla h$ ,  $h \in \mathbb{R}$ . It can be easily shown that Theorem 3.2.2 is equivalent to the preceding formulation in conjunction with  $h(w) = 0$ .

More often than not, one will encounter optimization problems with inequality constraints mixed in such as during a SVM problem. In light of this, we will introduce the generalized Lagrangian.

**Definition 3.2.2.** Given an optimization with domain  $\Omega \subseteq \mathbb{R}^n$ ,

$$\begin{aligned} \min_{w \in \Omega} \quad & f(w) \\ \text{s.t.} \quad & g_i(w) \leq 0, \quad i = 1, \dots, k, \\ & h_i(w) = 0, \quad i = 1, \dots, m, \end{aligned} \tag{3.2}$$

we define the *generalised Lagrangian function* as

$$L(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^m \beta_i h_i(w).$$

**Definition 3.2.3.** The *Lagrangian dual problem* of the primal problem of Definition 3.1.1 is the following problem:

$$\begin{aligned} \max \quad & \theta(\alpha, \beta) \\ \text{s.t.} \quad & \alpha \geq 0 \end{aligned} \tag{3.3}$$

where  $\theta(\alpha, \beta) = \inf_{w \in \Omega} L(w, \alpha, \beta)$ .

The relationship between the dual and primal problems will be clear after the following theorems.

**Theorem 3.2.3.** (*Weak duality theorem*) Let  $w \in \Omega$  be a feasible solution of the primal problem of 3.1.1 and  $(\alpha, \beta)$  a feasible solution of the dual problem of 3.3. Then  $f(w) \geq \theta(\alpha, \beta)$ .

*Proof.* By definition,

$$\begin{aligned} \theta(\alpha, \beta) &= \inf_{u \in \Omega} L(u, \alpha, \beta) \\ &\leq L(w, \alpha, \beta) \\ &= f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^m \beta_i h_i(w) \leq f(w), \end{aligned}$$

since  $g_i(w) \leq 0$  and  $h_i(w) = 0$ ,  $\alpha \geq 0$  by the feasibility of  $w$  and  $(\alpha, \beta)$ . □

Hence, it is not difficult to see that the value of the dual formation is upper bounded by the value of the primal,

$$\sup\{\theta(\alpha, \beta) : \alpha \geq 0\} \leq \inf\{f(w) : g(w) \leq 0, h(w) = 0\}.$$

**Corollary 3.2.1.** If  $f(w^*) = \theta(\alpha^*, \beta^*)$ , where  $(\alpha^*, \beta^*), w^*$  are feasible then they solve the dual and primal problems respectively. In this case,  $\alpha_i^* g_i(w^*) = 0$ , for  $i = 1, \dots, k$ .

*Proof.* Since  $f(w^*) = \theta(\alpha^*, \beta^*)$  and  $\alpha_i^* g_i(w^*) = 0$ , for  $i = 1, \dots, k$ ,

$$\begin{aligned} \theta(\alpha^*, \beta^*) &= f(w^*) \\ &= f(w^*) + \sum_{i=1}^k \alpha_i^* g_i(w^*) + \sum_{i=1}^m \beta_i^* h_i(w^*) \\ &= L(w^*, \alpha^*, \beta^*) \\ &= \inf_{u \in \Omega} L(u, \alpha^*, \beta^*). \end{aligned}$$

□

By comparing the primal and dual values in parallel and check the *duality gap* i.e their difference, one may be able to find the optimal solution if this reduces to zero. However, it is not generally guaranteed that the primal and dual problems will have the same values as solution.

**Theorem 3.2.4.** (*Strong duality theorem*) Given an optimization problem with convex domain  $\Omega \subset \mathbb{R}^n$ ,

$$\begin{aligned} \min_{w \in \Omega} \quad & f(w) \\ \text{s.t.} \quad & g_i(w) \leq 0, \quad i = 1, \dots, k, \\ & h_i(w) = 0, \quad i = 1, \dots, m, \end{aligned} \tag{3.4}$$

where  $g_i, h_i$  are affine functions and  $f$  convex then the duality gap is zero.

The following theorem given by Kuhn-Tucker states the conditions for an optimal solution to a general optimization problem.

**Theorem 3.2.5.** (*Kuhn-Tucker*) Given an optimization problem with convex domain  $\Omega \subseteq \mathbb{R}^n$ ,

$$\begin{aligned} \min_{w \in \Omega} \quad & f(w) \\ \text{s.t.} \quad & g_i(w) \leq 0, \quad i = 1, \dots, k, \\ & h_i(w) = 0, \quad i = 1, \dots, m, \end{aligned} \tag{3.5}$$

with  $f \in C^1$  convex and  $g_i, h_i$  affine. The necessary and sufficient conditions for a normal point  $w^*$  to be an optimum are the existence of  $\alpha^*, \beta^*$  such that

$$\begin{aligned} \frac{\partial L(w^*, \alpha^*, \beta^*)}{\partial w} &= 0, \\ \frac{\partial L(w^*, \alpha^*, \beta^*)}{\partial \beta} &= 0, \\ g_i(w^*) &\leq 0, \quad i = 1, \dots, k, \\ \alpha_i^* g_i(w^*) &= 0, \quad i = 1, \dots, k, \\ \alpha_i^* &\geq 0, \quad i = 1, \dots, k. \end{aligned}$$

The second and third conditions are necessary to satisfy the primal feasibility while the last condition satisfies the dual's. The fourth condition is usually known as Karush-Kuhn-Tucker complementary condition, it implies that for active constraints,  $\alpha_i^* \geq 0$ , whereas for inactive constraints  $\alpha_i^* = 0$ . This also follow from 3.2.1 since we are assuming zero duality gap.

The dual representation of a primal problem often turns out to be easier to solve since handling inequality constraints directly is difficult. The primal can be transformed into a dual by setting the derivatives of the Lagrangian w.r.t the primal variables and substituting the obtained relations back into the Lagrangian therefore removing the dependence on said variables. This corresponds to computing the function

$$\theta(\alpha, \beta) = \inf_{w \in \Omega} L(w, \alpha, \beta).$$

Which leaves us with the step of maximizing the resulting function under a (much) simpler constraint.

## Support Vector Machine

Given a set of (linearly separable) data points and labels like in Figure 1 . There are infinitely many hyperplanes that can be used to separate the two classes but optimally, we would want a hyperplane that can generalize well for unseen data. If a hyperplane that is too close to the data points is chosen then the permitted margin of error would be too small to predict unknown data accurately because it is easier for them to fall on either sides of the hyperplanes. Thus, an ideal separator would be the one that has the largest margin i.e the distance between the nearest data points to the plane is greatest [3].

### 4.1 Hard Margin SVM

We start off with the assumption that there exists a hyperplane that can perfectly separate or classifies the dataset.

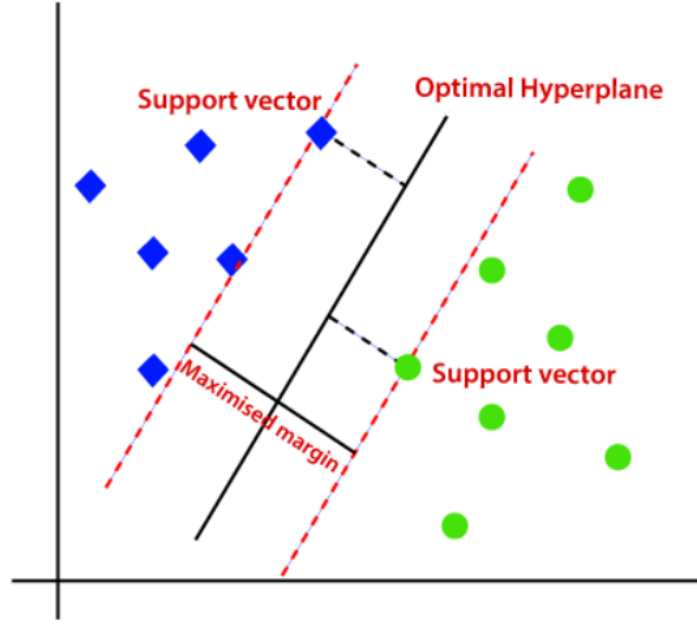


Figure 1: Support vectors [4]

To put it concretely, let  $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$  be the set of data points with  $x_i \in \mathbb{R}^m$  and  $y_i \in \{1, -1\}$  as the labels. Let  $h = w \cdot x + b$  be an arbitrary hyperplane then define the positive and negative support vectors,  $x^+$  and  $x^-$ , of this hyperplane as the points closest to  $h$ . Note that  $h, ch$ , for  $c \in \mathbb{R}$ , define the same hyperplane so we can choose our normalization factor such that

$$\begin{aligned}\langle w, x^+ \rangle + b &= 1 \\ \langle w, x^- \rangle + b &= -1.\end{aligned}$$

Then the set  $S$  is called linear separable if there exists a hyperplane defined by  $(w, b)$  such that

$$\begin{cases} \langle w, x_i \rangle + b \geq 1 & \text{for } y_i = 1 \\ \langle w, x_i \rangle + b \leq -1 & \text{for } y_i = -1 \end{cases}$$

or

$$y_i(\langle w, x_i \rangle + b) \geq 1 \quad \forall i \in \{1, 2, \dots, n\}.$$

The (geometric) margin is then given by

$$\gamma = \frac{|w \cdot x^+ + b|}{\|w\|} + \frac{|w \cdot x^- + b|}{\|w\|} = \frac{2}{\|w\|}$$

where  $\|\cdot\|$  is the standard Euclidean norm. Recall that the objective is to find a hyperplane such that  $\gamma$  is maximized while being subjected to  $y_i(w \cdot x_i + b) \geq 1$  i.e

$$\begin{aligned} \max_{w, b} \quad & \gamma \\ \text{s.t.} \quad & y_i(\langle w, x_i \rangle + b) \geq 1, \\ & i = 1, 2, \dots, n \end{aligned}$$

which is equivalent to

$$\begin{aligned} \min_{w, b} \quad & \frac{\|w\|^2}{2} \\ \text{s.t.} \quad & y_i(\langle w, x_i \rangle + b) \geq 1, \\ & i = 1, 2, \dots, n. \end{aligned} \tag{4.1}$$

i.e a quadratic optimization problem. The formulations stated above belong to the Hard Margin SVM problem because the hyperplane  $(w, b)$  must classify each point correctly. In practice, more often than not, this

is unachievable and impractical for a variety of reasons. Thus, ideally, it is desirable to have a classifier that allow for "minor" mistakes while retaining its generalization. Note that there are many versions to the Hard Margin SVM problem [5].

Following the steps outlined in Section 3, the primal Lagrangian is

$$L(w, b, \alpha) = \frac{\|w\|^2}{2} - \sum_{i=1}^n \alpha_i [y_i (\langle w, x_i \rangle + b) - 1]$$

where  $\alpha_i \geq 0$  are the Lagrange multipliers.

Differentiating the Lagrangian w.r.t to  $w, b$  to obtain

$$\begin{aligned} \frac{\partial L(w, b, \alpha)}{\partial w} &= w - \sum_{i=1}^n \alpha_i y_i x_i = 0 \implies w = \sum_{i=1}^n \alpha_i y_i x_i, \\ \frac{\partial L(w, b, \alpha)}{\partial b} &= \sum_{i=1}^n y_i \alpha_i = 0. \end{aligned}$$

Resubstituting the above relations into the primal Lagrangian to obtain

$$\begin{aligned} L(w, b, \alpha) &= \frac{\langle w, w \rangle}{2} - \sum_{i=1}^n \alpha_i [y_i (\langle w, x_i \rangle + b) - 1] \\ &= \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle - \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle + \sum_{i=1}^n \alpha_i \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle. \end{aligned}$$

We have the resulting dual problem

c

**Remark.** The relation  $w = \sum_{i=1}^n \alpha_i y_i x_i$  shows that the weights  $w$  can be described as a linear combination of the training data. This is called the dual representation.

**Remark.** Recall the Kuhn-Tucker complimentary condition 3.2.5 from the previous chapter that for a solution  $\alpha^*, (w^*, b^*)$  to be optimal, it must satisfy

$$\alpha_i^* [y_i (\langle w^*, x_i \rangle + b^*) - 1] = 0, \quad i = 1, \dots, n.$$

Define  $\langle w, x_i \rangle + b$  as the functional margin of  $x_i$  w.r.t  $(w, b)$ . This implies that for points  $x_i$  which the functional margin is 1 and lies closest to the hyperplane has non-zero  $\alpha_i^*$  and vice versa. Thus, the hyperplane is only determined by the points closest to it i.e, support vectors.

Having solved the dual problem to obtain the Lagrange multipliers, we can easily find the weights  $w$  but  $b$  is still unknown. In [6],  $b$  is calculated by taking the average of  $y_i - \langle w, x_i \rangle$  i.e

$$b = \frac{1}{S} \sum_{i=1}^S (y_i - \langle w, x_i \rangle)$$

where  $S$  is the number of support vectors. Other authors such as [7] only take the average of the nearest positive and negative support vectors

$$b = -\frac{\max_{y_i=-1} \langle w, x_i \rangle + \min_{y_i=1} \langle w, x_i \rangle}{2}.$$

## 4.2 Soft Margin SVM

Supposed that the data is not linearly separable, then the optimization problem 4.1 can not be solved since the condition  $y_i (\langle w, x_i \rangle + b)$  is not satisfiable. Also, if the data is linearly separable but the separating hyperplane does not leave a wide enough margin to account for unseen data then this would indeed be a bad classifier. Thus arise the need to alter the original Hard Margin SVM problem to account for the trade-off between errors and generalizability.

Instead of forcing our classifier to be correct at every point  $x_i$ , we introduce the *slack* variables  $\xi_i$  to our constraint as

$$\begin{aligned} \min_{w, b} \quad & \frac{1}{2} \langle w, w \rangle \\ \text{s.t.} \quad & y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i, \\ & \xi_i \geq 0, \\ & i = 1, 2, \dots, n. \end{aligned} \tag{4.2}$$

Clearly, we could let  $\xi_i \rightarrow \infty$  and the constraint would be satisfied but this has no added-value to the problem so we would want to constraint the values of  $\xi_i$ . There are many formulation of this approach but we shall use the one introduced by [7],

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} \langle w, w \rangle + C \sum_{i=1}^n \xi_i^2 \\ \text{s.t.} \quad & y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i, \\ & \xi_i \geq 0, \\ & i = 1, 2, \dots, n \end{aligned} \tag{4.3}$$

for some  $C \in \mathbb{R}$ . This problem is called  $L_2$ -norm soft margin. The Lagrangian and dual problem can then be obtained in the same way as the hard margin one.

The constant  $C$  was introduced as a way to control the importance of the *slack* variables. If  $C \rightarrow \infty$  then it is necessary that  $\xi_i \rightarrow 0$  for  $i = 1, \dots, n$  and thus become a Hard Margin SVM problem. Each different  $C$  can lead to a different classifier that may or may not be satisfactory so to find the optimal  $C$ , we have to try different values. Some of the recommended approaches can be found in [8] but they are beyond the scope of this paper.

So far, we have only assumed that the data is either linearly separable or linear separable but with a few outliers that can still be somewhat accurately classified using a separating hyperplane. What if, the data is non-linearly separable in such a way that a linear separator would not yield satisfactory results. Hence, it is necessary that we introduce a very important method to transform our feature vectors in such a way that they become linearly separable.

## Positive Definite Kernels and the Reproducing Kernel Hilbert Space

This section aims to introduce the *kernel* method that is used to transform a feature vector into a higher or infinite dimensional space with the aim of making the dataset linearly separable. We will also present related and necessary objects to supplement the kernel method such as Hilbert space, Reproducing Kernel Hilbert Space (RKHS) and the Representer Theorem. Although there are many publications that omit these definitions, we decline to do the same since they are crucial for the goal of this paper, that is Laplacian SVM.

### 5.1 Positive Definite Kernels

We begin with a few necessary definitions of basic spaces to build up to the Hilbert space.

**Definition 5.1.1.** Let  $X$  be a non-empty set with a distance function (or metric)  $d : X \times X \rightarrow \mathbb{R}^+$  then  $X$  is a metric space if for all  $x, y, z \in X$ ,  $d$  satisfies the following conditions:

1.  $d(x, y) = 0 \iff x = y$ .
2.  $d(x, y) = d(y, x)$ .
3.  $d(x, y) + d(y, z) \geq d(x, z)$ .

**Example 5.1.1.** The set of real numbers  $\mathbb{R}$  forms a metric space with  $d(x, y) = |x - y|$ .

**Definition 5.1.2.** Let  $V$  be a vector space over  $\mathbb{R}$ . A norm function on  $V$  is a function  $\|\cdot\| : V \rightarrow \mathbb{R}$  such that for all  $v, w$  in  $V$  and  $a \in \mathbb{R}$ ,

1.  $\|v\| \geq 0$  with  $\|v\| = 0 \iff v = 0$ .
2.  $\|av\| = |a| \|v\|$ .
3.  $\|v + w\| \leq \|v\| + \|w\|$ .



This norm induces a metric  $d(v, w) = \|v - w\|$  on  $V$ .

**Example 5.1.2.** The space  $\mathbb{R}^n$  is a vector space over  $\mathbb{R}$  with the  $L_2$  norm:

$$\|v\|_2 = \sqrt{v_1^2 + v_2^2 + \cdots + v_n^2}.$$

**Definition 5.1.3.** A normed vector space  $V$  is complete if every Cauchy sequence in  $V$  converges to a vector  $v \in V$  with respect to the norm.

**Example 5.1.3.** The normed vector space  $(\mathbb{Q}, \|\cdot\|)$  is not complete because if a sequence is recursively defined as

$$x_0 = \frac{4}{3}$$

$$\forall n \in \mathbb{N} : x_{n+1} = \frac{4 + 3x_n}{3 + 2x_n}$$

then it can be shown [9] that this sequence is Cauchy and converges to  $\sqrt{2}$  which is not in  $\mathbb{Q}$ .

**Definition 5.1.4.** Let  $V$  be a vector space over  $\mathbb{R}$ . An inner product on  $V$  is a function  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$  such that for any  $v, w, u \in V$  and  $a, b \in \mathbb{R}$ , we have

1.  $\langle u, v \rangle = \langle v, u \rangle$ .
2.  $\langle au + bv, w \rangle = a\langle u, w \rangle + b\langle v, w \rangle$ .
3.  $\langle v, v \rangle \geq 0$  with  $\langle v, v \rangle = 0 \iff v = 0$ .

A vector space that is equipped with an inner product is called an inner product space.

**Example 5.1.4.** The vector space  $\mathbb{R}^n$  is an inner product space with the inner product as the dot product.

**Definition 5.1.5.** Let  $(V, \langle \cdot, \cdot \rangle)$  be an inner product space, the norm induced by  $\langle \cdot, \cdot \rangle$  can be defined as

$$\|v\| = \sqrt{\langle v, v \rangle} \text{ for } v \in V.$$

**Remark.** Not every norm is induced by an inner product. In fact, it holds if and only if the parallelogram identity

$$\|x + y\|^2 + \|x - y\|^2 = 2\|x\|^2 + 2\|y\|^2$$

holds.

**Definition 5.1.6.** (Hilbert space) A Hilbert space is an inner product space that is complete with respect to the norm induced by the inner product.

Hilbert space is important to our research because of its rich structure and the methods that we are going to introduce, all operate within Hilbert space. We now can define what is a kernel.

**Definition 5.1.7.** (Kernel)

**Example 5.1.5.** Let  $X \subset \mathbb{R}^2$ , and let  $\varphi : X \rightarrow \mathbb{R}^3$  be the map given by

$$\varphi(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2}x_1x_2).$$

We have

$$\begin{aligned} \langle \varphi(x_1, x_2), \varphi(y_1, y_2) \rangle &= \langle (x_1^2, x_2^2, \sqrt{2}x_1x_2), (y_1^2, y_2^2, \sqrt{2}y_1y_2) \rangle \\ &= x_1^2y_1^2 + x_2^2y_2^2 + 2x_1x_2y_1y_2 \\ &= (x_1y_1 + x_2y_2)^2 = \langle x, y \rangle, \end{aligned}$$

where  $\langle x, y \rangle$  is the usual inner product on  $\mathbb{R}^2$ . Hence,

$$K(x, y) = \langle x, y \rangle^2$$

is a kernel function associated with the feature space  $\mathbb{R}^3$ .

One of the reason why the kernel function is often used in various contexts of Machine Learning is that, as you have seen, the calculation and derivation of the feature map  $\varphi$  is time-consuming and resource intensive while it could be replaced by a simple inner product. Furthermore, if the Hilbert space  $H$  is infinite dimensional then clearly, computing  $\varphi$  is impossible but nevertheless, the expression  $K(x, y)$  can be found, avoiding  $\varphi$  altogether.

Using the kernel function, we can then transform the hard margin dual problem into

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j K(x_i, x_j) \\ \text{s.t.} \quad & \alpha_i \geq 1 \quad i = 1, \dots, n, \\ & \sum_{i=1}^n y_i \alpha_i = 0 \quad i = 1, \dots, n \end{aligned} \tag{5.1}$$

which can then be solved similarly to the old one with the exception that the dataset is now linearly separable, assuming that it is indeed the case in the feature space.

**Remark.** The same kernel can arise from different maps into different feature spaces. Example:  $\varphi_1(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$  in  $\mathbb{R}^3$  and  $\varphi_2(x_1, x_2) = (x_1^2, x_2^2, x_1x_2, x_1x_2)$  in  $\mathbb{R}^4$ .

One of the important property of the kernel functions is the positive definite-ness and the reason for its importance will become clear by the end of this section.

**Definition 5.1.8.** (Positive definite kernel) Let  $X$  be a non-empty set. A function  $\kappa : X \times X \rightarrow \mathbb{C}$  is a *positive definite kernel* if for every finite  $S \subset X$ , if  $K_s$  is the  $p \times p$  matrix

$$K_s = (K(x_j, x_i))_{1 \leq i, j \leq p}$$

then we have

$$u^* K_s u = \sum_{i,j=1}^p K(x_i, x_j) u_i \overline{u_j} \geq 0, \quad \text{for all } u \in \mathbb{C}^p$$

where  $u^*$  denotes the complex conjugate of  $u$ .

**Remark.** A positive definite kernel that is also symmetric i.e  $K(x, y) = K(y, x)$  is called a *Mercer Kernel*.

**Proposition 5.1.1.** Let  $X$  be any non-empty set, let  $H$  be any (complex) Hilbert space, let  $\varphi : X \rightarrow H$  be any function, and let  $K : X \times X \rightarrow \mathbb{C}$  be the kernel given by

$$K(x, y) = \langle \varphi(x), \varphi(y) \rangle, \quad x, y \in X.$$

For any finite subset  $S = \{x_1, \dots, x_p\}$  of  $X$ , if  $K_s$  is the  $p \times p$  matrix

$$K_s = (K(x_j, x_i))_{1 \leq i, j \leq p}$$

then we have

$$u^* K_s u \geq 0, \quad \text{for all } u \in \mathbb{C}^p.$$

*Proof.*

$$\begin{aligned} u^* K_s u &= \sum_{i,j=1}^p K(x_i, x_j) u_i \overline{u_j} \\ &= \sum_{i,j=1}^p \langle \varphi(x_i), \varphi(x_j) \rangle u_i \overline{u_j} \\ &= \left\langle \sum_{i=1}^p u_i \varphi(x_i), \sum_{j=1}^p u_j \varphi(x_j) \right\rangle \\ &= \left\| \sum_{i=1}^p u_i \varphi(x_i) \right\|^2 \geq 0. \end{aligned}$$

□

**Proposition 5.1.2.** (I. Schur) If  $K_1, K_2 : X \times X \rightarrow \mathbb{C}$  are two positive definite kernels, then the function  $K : X \times X \rightarrow \mathbb{C}$  given by  $K(x, y) = K_1(x, y) K_2(x, y)$  for all  $x, y \in X$  is also a positive definite kernel.

*Proof.* [5]. □

Here are some ways of obtaining new positive definite kernels from old ones

**Proposition 5.1.3.** *Let  $K_1 : X \times X \rightarrow \mathbb{C}$  and  $K_2 : X \times X \rightarrow \mathbb{C}$  be two positive definite kernels, let  $f : X \rightarrow \mathbb{C}$ ,  $\psi : X \rightarrow \mathbb{R}^N$  be functions,  $K_3 : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{C}$  be a positive definite kernel, and  $a \in \mathbb{R}^+$ ,  $p(z)$  be a polynomial with nonnegative coefficients. Then the following functions are positive definite kernels:*

- (1)  $K(x, y) = K_1(x, y) + K_2(x, y)$ .
- (2)  $K(x, y) = aK_1(x, y)$ .
- (3)  $K(x, y) = K_3(\psi(x), \psi(y))$ .
- (4)  $K(x, y) = f(x)\overline{f(y)}$ .
- (5)  $K(x, y) = p(K_1(x, y))$ .
- (6)  $K(x, y) = e^{K_1(x, y)}$ .
- (7) If  $X$  is a real Hilbert space with inner product  $\langle -, - \rangle_X$  and corresponding norm  $\| \cdot \|_X$ ,

$$K(x, y) = e^{-\frac{\|x-y\|_X^2}{2\sigma^2}}$$

for any  $\sigma > 0$ .

*Proof.* (1), (2) and (3) are trivial.

- (4) Let  $S = \{x_1, \dots, x_p\} \subset X$ , if  $K$  is the  $p \times p$  matrix

$$K = (\overline{f(x_k)}f(x_j))_{1 \leq j, k \leq p}$$

then we have

$$u^*Ku = \sum_{k,j=1}^p u_j f(x_j) \overline{u_k f(x_k)} = \left| \sum_{j=1}^p u_j f(x_j) \right|^2 \geq 0.$$

- (5) Let  $p(z) = \sum_{i=0}^m a_i z^i$ , then

$$p(K_1(x, y)) = a_m K_1(x, y)^m + \dots + a_1 K_1(x, y) + a_0.$$

Since  $a_i \in \mathbb{R}^+$  for  $i = 0, \dots, m$ , by Proposition 5.1.2 and (2), each  $a_i K_1(x, y)^i$  is a positive definite kernel. By (4), (1) with  $f(x) = \sqrt{a_0}$ ,  $p(K_1(x, y))$  is a positive definite kernel.

- (6) We first show that if each  $K_i : X \times X \rightarrow \mathbb{C}$  is a positive definite kernel,

$$\lim_{i \rightarrow \infty} K_i(x, y) = K(x, y)$$

is also a positive definite kernel if it exists. We have

$$\begin{aligned} \sum_{j,k=1}^n u_j u_k K(x_j, x_k) &= \sum_{j,k=1}^n u_j u_k \left( \lim_{i \rightarrow \infty} K_i(x_j, x_k) \right) \\ &= \sum_{j,k=1}^n \lim_{i \rightarrow \infty} (u_j u_k K_i(x_j, x_k)) \\ &= \lim_{i \rightarrow \infty} \underbrace{\left( \sum_{j,k=1}^n u_j u_k K_i(x_j, x_k) \right)}_{\geq 0} \geq 0. \end{aligned}$$

Thus,  $K$  is positive definite. Note that

$$e^{K_1(x, y)} = \lim_{n \rightarrow \infty} \sum_{k=0}^n \frac{K_1(x, y)^k}{k!}$$

but each partial sum

$$\sum_{k=0}^n \frac{K_1(x, y)^k}{k!}$$

is a positive definite kernel so  $e^{K_1(x, y)}$  is also a positive definite kernel.

(7) By (2) and since the map  $(x, y) \rightarrow \langle x, y \rangle_X$  is a positive definite kernel by Proposition 5.1.1 with identity feature map, the function

$$(x, y) \rightarrow \frac{\langle x, y \rangle_X}{\sigma^2}$$

is a positive definite kernel so it follows that

$$K_1 = e^{\frac{\langle x, y \rangle_X}{\sigma^2}}$$

is also a positive definite kernel. Let  $f : X \rightarrow \mathbb{R}$  be defined as

$$f(x) = e^{-\frac{\|x\|_X^2}{2\sigma^2}}$$

then by (4),

$$K_2(x, y) = f(x)\overline{f(y)} = f(x)f(y) = e^{-\frac{\|x\|_X^2}{2\sigma^2}} e^{-\frac{\|y\|_X^2}{2\sigma^2}} = e^{-\frac{\|x\|_X^2 + \|y\|_X^2}{2\sigma^2}}$$

is a positive definite kernel. Thus,

$$\begin{aligned} K_1(x, y)K_2(x, y) &= e^{\frac{2\langle x, y \rangle_X}{2\sigma^2}} e^{-\frac{\|x\|_X^2 + \|y\|_X^2}{2\sigma^2}} \\ &= e^{-\frac{\|x - y\|_X^2}{2\sigma^2}} \end{aligned}$$

is a positive definite kernel. It is usually called *Gaussian kernel*. □

## 5.2 Reproducing Kernel Hilbert Space

We are ready to introduce the Reproducing Kernel Hilbert Space (RKHS) which is a Hilbert space of functions with the reproducing property. The reason why this space is important has to do with the *Moore-Aronszajn theorem* which connect positive definite kernels to feature space and allow the kernel trick to work.

**Definition 5.2.1.** (RKHS) A Reproducing Kernel Hilbert Space is a Hilbert space  $\mathcal{H}$  of functions  $f : X \rightarrow \mathbb{R}$  with a reproducing kernel  $K : X \times X \rightarrow \mathbb{R}$  where  $K(x, \cdot) \in \mathcal{H}$  and  $f(x) = \langle K(x, \cdot), f \rangle_{\mathcal{H}}$  [10].

The property  $f(x) = \langle K(x, \cdot), f \rangle$  is called the *reproducing property*.

**Theorem 5.2.1** (Moore-Aronszajn). *If  $K$  is a Mercer kernel on a set  $X$ . Then there is a unique RKHS space of functions on  $X$  for which  $K$  is the reproducing kernel.*

*Proof.* We will only provide a sketch of the proof. The full details can be found at [11].

Let  $x \in X$  and define  $K_x = K(x, \cdot)$ . Let  $H_0$  be the linear span of  $\{K_x : x \in X\}$ . Let  $f, g \in H_0$  given by

$$f = \sum_{i=1}^m a_i K_{x_i}, \quad g = \sum_{j=1}^n b_j K_{y_j},$$

then define an inner product on  $H_0$  as

$$\langle f, g \rangle_{H_0} = \sum_{i=1}^m \sum_{j=1}^n a_i b_j K(x_i, y_j).$$

This shows that the inner product is symmetric bilinear. Using Cauchy-Schwartz inequality, it can be shown that  $H_0$  is a pre-Hilbert space. Let  $H$  be the set of functions which are pointwise limits of Cauchy sequences in  $H_0$ . We can then show that  $H$  is complete w.r.t the inner product on  $H_0$  then  $H$  is a RKHS which admit  $K$  as its reproducing kernel and every  $f \in H$  can be written as a linear combination of  $K$ . □

This theorem combines with the definition of RKHS states that for each RKHS we have a reproducing Mercer kernel and for each Mercer kernel, we can get a unique RKHS.

To put it more concretely, suppose that instead of wanting to compute  $\langle x_i, x_j \rangle$  in the dual formulation of the Hard Margin SVM problem

$$L(w, b, a) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle,$$

we want to compute the inner product of the feature map  $\varphi : X \rightarrow H$ , for  $H$  a Hilbert space,  $\langle \varphi(x_i), \varphi(x_j) \rangle$  in a higher dimensional space. This is potentially problematic because it requires us to explicitly compute the feature map which can be computationally expensive or impossible if the dimension of  $H$  is infinite. However, if instead one chooses a Mercer kernel  $K$  then by 5.2.1, there exists a unique RKHS  $\mathcal{H}$  where  $K$  has the reproducing property. Then this becomes much easier by for example, trivially letting  $\varphi(x) = K_x \in \mathcal{H}$  then

$$K(x_i, x_j) = \langle K_{x_i}(\cdot), K_{x_j}(\cdot) \rangle = \langle \varphi(x_i), \varphi(x_j) \rangle_{\mathcal{H}},$$

thus skipping the step of computing the feature map. However, Theorem 5.2.1 only guarantees the existence of feature maps but not how to find them, and for that we shall require one important theorem that provides a representation of the feature map.

Note that for a given Mercer kernel  $K$ , the corresponding feature map is not unique.

**Example 5.2.1.** Let  $K$  be a Mercer kernel given by  $K(x, y) = \langle x, y \rangle^2$  with  $X = \mathbb{R}^2$  then

$$K(x, y) = x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 x_2 y_1 y_2.$$

It can be easily seen that the two feature maps

$$\begin{aligned} \phi(x) &= (x_1^2, x_2^2, \sqrt{2}x_1 x_2) \\ \varphi(x) &= (x_1^2, x_2^2, x_1 x_2, x_1 x_2) \end{aligned}$$

correspond to the same kernel  $K$  with different feature spaces.

**Definition 5.2.2.** ( $L_p$  space) Consider a function  $f$  with domain  $[a, b] \subset \mathbb{R}$ . For  $p > 0$ , let the  $L_p$  norm be defined as :

$$\|f\|_p = \left( \int |f(x)|^p dx \right)^{\frac{1}{p}}.$$

The  $L_p$  space is defined as the set of functions with bounded  $L_p$  norm:

$$L_p(a, b) := \{f : [a, b] \rightarrow \mathbb{R} \mid \|f\|_p < \infty\}.$$

**Theorem 5.2.2** (Mercer's Theorem). Suppose  $K : [a, b] \times [a, b] \rightarrow \mathbb{R}$  is a continuous symmetric positive semi-definite kernel which is bounded:

$$\sup_{x, y} K(x, y) < \infty.$$

Define the operator  $T_k$  as

$$T_k f(x) = \int_a^b K(x, y) f(y) dy.$$

The operator  $T_k$  is called the Hilbert-Schmidt integral operator. This output function is positive definite:

$$\iint K(x, y) f(y) dx dy \geq 0.$$

Then there is a set of orthonormal basis  $\{\phi_i(\cdot)\}_{i=0}^{\infty}$  of  $L_2(a, b)$  consisting of eigenfunctions of  $T_k$  such that the corresponding sequence of eigenvalues  $\{\lambda_i\}_{i=1}^{\infty}$  are non-negative:

$$\int K(x, y) \phi_i(y) dy = \lambda_i \phi_i(x).$$

The eigenfunctions corresponding to the non-zero eigenvalues are continuous on  $[a, b]$  and  $k$  can be represented as

$$K(x, y) = \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(y),$$

where the convergence is absolute and uniform.

One may notice the similarity between Mercer and Moore-Aronszajn theorems that they both show there exist feature maps for a Mercer kernel  $K$ . The difference is Mercer's theorem creates a feature map by computing the eigenfunctions of  $T_k$  while Moore-Aronszajn's let the feature map  $\phi = K_x$ .

## Riemannian Manifold

A manifold is a mathematical structure that locally resembles Euclidean space. This structure is useful in the field of Machine Learning is partly because of the Manifold Hypothesis: many high-dimensional data sets that occur in the real world actually lie along low-dimensional manifolds inside that high-dimensional space. Manifold regularization is a technique using the shape of the data set to constraint on the function that should be learned on it. This technique applies a constraint on the Laplace-Beltrami operator, defined on a Riemannian Manifold, of functions in the RKHS of functions from an arbitrary set to the real numbers. We start by introducing the basic ideas of topological manifolds, smooth manifolds to its tangent vectors, and then Riemannian manifolds and the Laplace-Beltrami operator.

### 6.1 Topology prerequisites

**Definition 6.1.1.** A topology on a set  $X$  is a collection  $\mathcal{T}$  of subsets of  $X$  having the following properties:

- (1)  $\emptyset$  and  $X$  are in  $\mathcal{T}$ .
- (2) The union of the elements of any subcollection of  $\mathcal{T}$  is in  $\mathcal{T}$ .
- (3) The intersection of the elements of any finite subcollection of  $\mathcal{T}$  is in  $\mathcal{T}$ .

A set  $X$  for which a topology  $\mathcal{T}$  has been specified is called a topology space.

If  $X$  is a topological space with topology  $\mathcal{T}$ , we say that a subset  $U$  of  $X$  is an *open set* of  $X$  if  $U$  belongs to the collection  $\mathcal{T}$ .

**Definition 6.1.2.** If  $X$  is a set, a basis for a topology on  $X$  is a collection  $\mathcal{B}$  of subsets of  $X$  such that

- (1) For each  $x \in X$ , there is at least one basis element  $B$  containing  $x$ .
- (2) If  $x$  belongs to the intersection of two basis elements  $B_1$  and  $B_2$ , then there is a basis elements  $B_3$  containing  $x$  such that  $B_3 \subset B_1 \cap B_2$ .

If  $\mathcal{B}$  satisfies these two conditions, then we define the topology  $\mathcal{T}$  generated by  $\mathcal{B}$  as follows: a subset  $U$  is said to be open in  $X$  if for each  $x \in U$ , there is a basis element  $B \in \mathcal{B}$  such that  $x \in B$  and  $B \subset U$ . Note that each basis element is itself an element of  $\mathcal{T}$ .

**Definition 6.1.3.** Let  $X$  be a topological space with topology  $\mathcal{T}$ . If  $Y \subset X$ , the collection

$$\mathcal{T}_Y = \{Y \cap U \mid U \in \mathcal{T}\}.$$

is a topology on  $Y$ , called the subspace topology. With this topology,  $Y$  is called a subspace of  $X$ ; its open sets consist of all intersections of open sets of  $X$  with  $Y$ .

If  $Y$  is a subspace of  $X$ , we say that a set  $U$  is open in  $Y$  if it belongs to the topology of  $Y$ . We say that  $U$  is open in  $X$  if it belongs to the topology of  $X$ .

**Lemma 6.1.1.** Let  $Y$  be a subspace of  $X$ . If  $U$  is open in  $Y$  and  $Y$  is open in  $X$ , then  $U$  is open in  $X$ .

*Proof.* Since  $U$  is open in  $Y$ ,  $U = Y \cap V$  for some set  $V$  open in  $X$ . Since  $Y$  and  $V$  are both open, so is  $Y \cap V$ .  $\square$

**Definition 6.1.4.** (Hausdorff space)

A topological space  $X$  is called a Hausdorff space if for each pair  $x_1, x_2$  of distinct points of  $X$ , there exist  $U_1, U_2$  such that  $x_1 \in U_1$ ,  $x_2 \in U_2$  and  $U_1 \cap U_2 = \emptyset$ .

**Definition 6.1.5.** (Continuous function)

Let  $X, Y$  be topological spaces. A function  $f : X \rightarrow Y$  is said to be continuous if for each open subset  $V$  of  $Y$ , the set  $f^{-1}(V)$  is an open subset of  $X$ .

Continuity of a function depends not only upon the function  $f$  itself, but also on the topologies specified for its domain and range.

**Theorem 6.1.1.** *Let  $X, Y$ , and  $Z$  be topological spaces.*

- (a) *If  $A$  is a subspace of  $X$ , the inclusion function  $j : A \rightarrow X$  is continuous.*
- (b) *if  $f : X \rightarrow Y$  and  $g : Y \rightarrow Z$  are continuous, then the map  $g \circ f : X \rightarrow Z$  is continuous.*
- (c) *If  $f : X \rightarrow Y$  is continuous, and if  $A$  is a subspace of  $X$ , then the restricted function  $f|_A : A \rightarrow Y$  is continuous.*
- (d) *Let  $f : X \rightarrow Y$  be continuous. If  $Z$  is a subspace of  $Y$  containing  $f(X)$ , then the function  $g : X \rightarrow Z$  obtained by restricting the codomain of  $f$  is continuous.*

*Proof.* (a) If  $U$  is open in  $X$ , then  $j^{-1}(U) = U \cap A$ , which is open in  $A$  by the definition of the subspace topology.  
(b) If  $U$  is open in  $Z$ , then  $g^{-1}(U)$  is open in  $Y$  and  $f^{-1}(g^{-1}(U))$  is open in  $X$  but

$$f^{-1}(g^{-1}(U)) = (g \circ f)^{-1}(U).$$

(c) We have

$$f|_A = f \circ j$$

which is continuous.

(d) Let  $B$  be open in  $Z$ . Then  $B = Z \cap U$  for some open set  $U$  of  $Y$ . Because  $Z$  contains the entire image set  $f(X)$ ,

$$f^{-1}(U) = g^{-1}(B).$$

Since  $f^{-1}(U)$  is open, so is  $g^{-1}(B)$ . □

**Definition 6.1.6.** (Homeomorphism)

Let  $X$  and  $Y$  be topological spaces; let  $f : X \rightarrow Y$  be a bijection. If both the function  $f$  and the inverse function

$$f^{-1} : Y \rightarrow X$$

are continuous, then  $f$  is called a homeomorphism.

**Corollary 6.1.1.** Let  $f : X \rightarrow Y$  be a homeomorphism and  $A$  an open subset of  $X$ . The restriction of  $f|_A$  is a homeomorphism.

*Proof.* This is immediate from 6.1.1 (c) and (d). □

## 6.2 Smooth manifolds

**Definition 6.2.1.** (Topological manifold) Let  $M$  be a topological space. We say that  $M$  is a topological manifold of dimension  $n$  or a topological  $n$ -manifold if it has the following properties:

- $M$  is a Hausdorff space.
- $M$  is second-countable: there exists a countable basis for the topology of  $M$ .
- $M$  is locally Euclidean of dimension  $n$ : each point of  $M$  has a neighborhood that is homeomorphic to an open subset of  $\mathbb{R}^n$ .

The third property means, more specifically, that for each  $p \in M$  we can find

- an open subset  $U \subseteq M$  containing  $p$ ,
- an open subset  $\hat{U} \subseteq \mathbb{R}^n$  and
- a homeomorphism  $\varphi : U \rightarrow \hat{U}$ .

Every topological manifold has, by definition, a specific well-defined dimension.

**Theorem 6.2.1.** *A nonempty  $n$ -dimensional topological manifold cannot be homeomorphic to an  $m$ -dimensional manifold unless  $m = n$  [12].*

The basic example of a topological  $n$ -manifold is  $\mathbb{R}^n$  itself with the identity as homeomorphism. The reasons why manifolds have to satisfy the three properties is to ensure that they behave in a way that is similar to Euclidean space.

**Proposition 6.2.1.** *Every open subset of a topological  $n$ -manifold is itself a topological  $n$ -manifold.*

*Proof.* Let  $M$  be a topological  $n$ -dimensional manifold and  $U$  an open subset of  $M$ . It suffices to show that  $U$  is locally Euclidean since the other two requirements are trivial. Let  $x$  be a point in  $U$  then there exists a neighborhood  $V$  containing  $x$  that is homeomorphic to an open subset  $\hat{V} \subseteq \mathbb{R}^n$ . W.l.o.g let  $\varphi : V \rightarrow \hat{V}$  be a homeomorphism such that  $\varphi(x) = 0$ . It is clear that  $U \cap V = N$  is an open subset of  $V$  so  $\varphi(N) \subseteq \mathbb{R}^n \ni 0$ . Let  $0 \in B \subset \varphi(N)$  then  $\varphi^{-1}(B)$  is a neighborhood of  $x$  contained in  $V$ , which is homeomorphic to  $B \subseteq \mathbb{R}^n$ .  $\square$

**Definition 6.2.2.** (Coordinate chart) Let  $M$  be a topological  $n$ -manifold. A coordinate chart on  $M$  is a pair  $(U, \varphi)$ , where  $U$  is an open subset of  $M$  and  $\varphi : U \rightarrow \hat{U}$  is a homeomorphism.

Thus far, we have yet to define any calculus operations on manifolds such as derivatives, etc which are the heart of machine learning. To do that, we will introduce a structure called *smooth manifold*. The definition will be based on maps between Euclidean spaces and from now on, we define a map between open subsets of Euclidean spaces to be smooth if all of its component functions has continuous partial derivatives of all orders.

**Definition 6.2.3.** (Diffeomorphism) Given two topological manifolds  $M$  and  $N$ . A smooth map  $f : M \rightarrow N$  is a *diffeomorphism* if it is bijective and its inverse is also smooth.

**Definition 6.2.4.** (Transition map) Let  $M$  be a topological  $n$ -manifold. If  $(U, \varphi)$ ,  $(V, \psi)$  are two charts such that  $U \cap V \neq \emptyset$ , the composite map  $\psi \circ \varphi^{-1} : \varphi(U \cap V) \rightarrow \psi(U \cap V)$  is called the *transition map* from  $\varphi$  to  $\psi$ . Two charts  $(U, \varphi)$  and  $(V, \psi)$  are said to be *smoothly compatible* if either  $U \cap V = \emptyset$  or the transition map is a diffeomorphism.

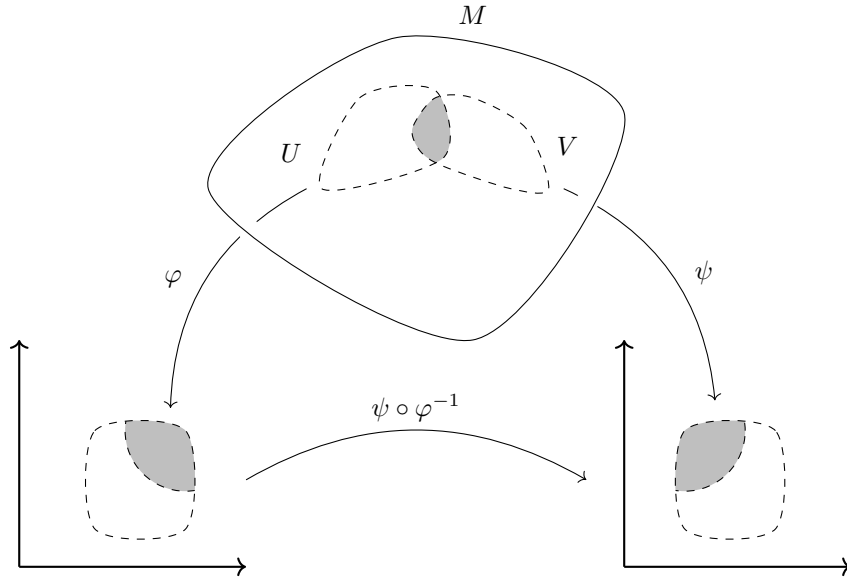


Figure 2: Transition map

Since the transition map has domain and codomain as open subsets of Euclidean spaces, its smoothness can be defined as having continuous partial derivatives of all orders.

**Definition 6.2.5.** (Smooth Atlas) An atlas for a topological manifold  $M$  is a collection of charts whose domain cover  $M$ . An atlas  $\mathcal{A}$  is called a *smooth atlas* if for any two charts in  $\mathcal{A}$ , they are smoothly compatible.

Note that for a topological manifold  $M$ , there may exist different atlases that determine the same collection of smooth functions on  $M$ . To that end, we shall define *maximal atlas*.

**Definition 6.2.6.** (Maximal Smooth Atlas) A smooth atlas  $\mathcal{A}$  on  $M$  is maximal if it is not properly contained in any larger smooth atlas.

**Definition 6.2.7.** (Smooth manifold) If  $M$  is a topological manifold, a *smooth structure* on  $M$  is a maximal smooth atlas. Then the pair  $(M, \mathcal{A})$  is called a *smooth manifold*.

**Example 6.2.1.** For each non-negative integer  $n$ , the space  $\mathbb{R}^n$  is a smooth  $n$ -manifold with its smooth structure determined by a single chart  $(\mathbb{R}^n, \text{Id}_{\mathbb{R}^n})$ . This is called the standard smooth structure on  $\mathbb{R}^n$ .

If  $M$  is a smooth manifold, any chart  $(U, \varphi)$  contained in the given maximal smooth atlas is called a *smooth chart*, and the corresponding coordinate map  $\varphi$  is called a *smooth coordinate map*. We can represent a point  $p \in U$  by its local coordinates  $\varphi(p) = (x^1, \dots, x^n)$  using Einstein notation. It is customary to say that  $p = (x^1, \dots, x^n)$  in local coordinates.



### 6.3 Smooth maps

From now on, any manifold is considered smooth unless stated otherwise. We usually use "functions" as maps from an object to an Euclidean space and "map" or "mapping" as maps between any arbitrary objects, in this case, manifolds.

**Definition 6.3.1.** (Smooth functions) Let  $M$  be a manifold,  $k$  a non-negative integer, and  $f : M \rightarrow \mathbb{R}^k$  any function. We say that  $f$  is a *smooth function* if for every  $p \in M$ , there exists a smooth chart  $(U, \varphi)$  for  $M$  whose domain contains  $p$  and such that the function  $f \circ \varphi^{-1}$  is smooth on  $\varphi(U) \subseteq \mathbb{R}^n$ .

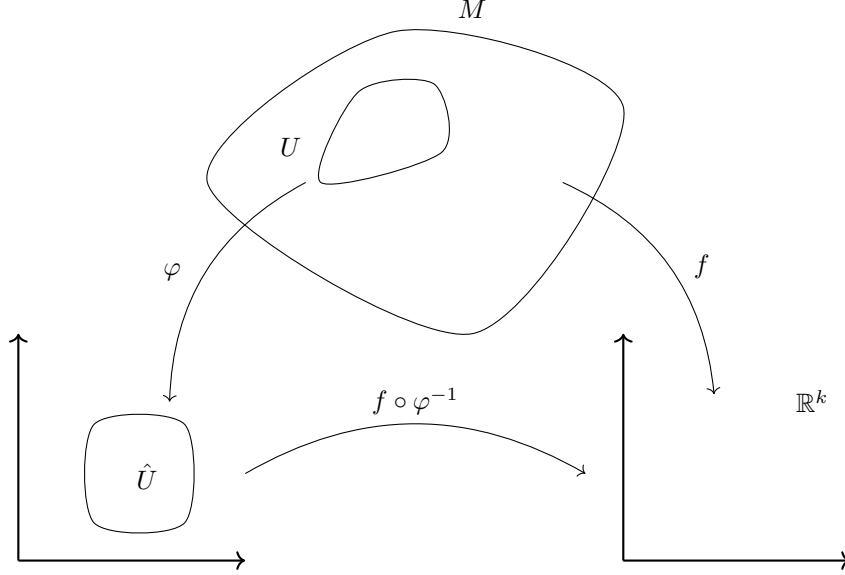


Figure 3: Smooth functions

The map  $\hat{f} : \varphi(U) \rightarrow \mathbb{R}^k$  defined by  $\hat{f} = f \circ \varphi^{-1}$  is called the *coordinate representation* of  $f$ . By definition,  $f$  is smooth if and only if its coordinate representation is smooth in some smooth chart around each point. We can easily generalize the definition of smooth functions to maps between manifolds.

**Definition 6.3.2.** (Smooth maps) Let  $M, N$  be manifolds, and let  $F : M \rightarrow N$  be any map. We say that  $F$  is a *smooth map* if for every  $p \in M$ , there exists a smooth chart  $(U, \varphi)$  containing  $p$  and  $(V, \psi)$  containing  $F(p)$  such that  $F(U) \subseteq V$  and the composite map  $\psi \circ F \circ \varphi^{-1}$  is smooth from  $\varphi(U)$  to  $\psi(V)$ .

### 6.4 Tangent Vectors

In this section, we will introduce *tangent space* to a manifold at a point to help us better define calculus on manifolds and *differential* which generalizes the total derivative of a map between Euclidean spaces. We start off with concrete objects like geometric tangent vectors in  $\mathbb{R}^n$  given by [12].

Given a point  $a \in \mathbb{R}^n$ , define the geometric tangent space to  $\mathbb{R}^n$  at  $a$ , denoted as  $\mathbb{R}_a^n$ , to be the set  $\{(a, v) : v \in \mathbb{R}^n\}$ . A geometric tangent vector in  $\mathbb{R}^n$  is an element of  $\mathbb{R}_a^n$  for some  $a \in \mathbb{R}^n$ . Clearly, the set  $\mathbb{R}_a^n$  is a vector space under standard operations

$$cv_a + cw_a = c(v + w)_a$$

where  $v_a$  is a vector starting at  $a$ . Thus, with this definition, one could think of the tangent space to a sphere with ambient space  $\mathbb{R}^3$  at a point  $a$  on the sphere as a space of vectors that are orthogonal to the radial unit vector through  $a$  using the inherited natural inner product from  $\mathbb{R}^n$ . The glaring problem with this definition is that in an arbitrary manifold with non-Euclidean ambient space, we have no idea how to take an inner product or how to define orthogonality.

Using the directional derivative of functions as an inspiration, we call a map  $w : C^\infty(\mathbb{R}^n) \rightarrow \mathbb{R}$  as a *derivation* at  $a$  if it is linear over  $\mathbb{R}$  and satisfies the following:

$$w(fg) = f(a)w(g) + g(a)w(f).$$

Let  $T_a\mathbb{R}^n$  denote the set of all derivations of  $C^\infty(\mathbb{R}^n)$  at  $a$ . It can be seen that  $T_a\mathbb{R}^n$  is a vector space under the operations

$$(w_1 + w_2)(f) = w_1(f) + w_2(f).$$

**Lemma 6.4.1.** Suppose  $a \in \mathbb{R}^n$ ,  $w \in T_a \mathbb{R}^n$ , and  $f, g \in C^\infty(\mathbb{R}^n)$ .

(a) If  $f$  is a constant function, then  $w(f) = 0$ .

(b) If  $f(a) = g(a) = 0$ , then  $w(fg) = 0$ .

*Proof.* (a) Let  $f_1(x) \equiv 1$  then  $f(x) \equiv cf_1$ . We have

$$w(f_1) = w(f_1 f_1) = f_1 w(f_1) + f_1 w(f_1) = 2w(f_1),$$

which implies that  $w(f_1) = 0$ . It follows that  $w(f) = w(cf_1) = cw(f_1) = 0$ .

(b) By the product rule,

$$w(fg) = f(a)w(g) + g(a)w(f) = 0.$$

□

**Proposition 6.4.1.** Let  $a \in \mathbb{R}^n$ .

(a) For each geometric tangent vector  $v_a \in \mathbb{R}_a^n$ , the map  $D_{v|_a} : C^\infty(\mathbb{R}^n) \rightarrow \mathbb{R}$  defined as

$$D_{v|_a} f = D_v f(a) = \left. \frac{d}{dt} \right|_{t=0} f(a + tv)$$

is a derivation at  $a$ .

(b) The map  $v_a \rightarrow D_{v|_a}$  is an isomorphism from  $\mathbb{R}_a^n$  onto  $T_a \mathbb{R}^n$ .

*Proof.* (a) This is immediate from the definition of derivative.

(b) Note that the map is linear. Suppose  $v_a \in \mathbb{R}_a^n$  such that  $D_{v|_a}$  is the zero derivation. Let  $v_a = v^i e_i|_a$  in Einstein notation ( $v^i e_i|_a := \sum_i v^i e_i|_a$ ) where  $x^i$  means the  $i$ -th component of  $x$  and  $e_i$  is the standard  $i$ -th basis. Let  $f$  be the  $x^j : \mathbb{R}^n \rightarrow \mathbb{R}$  smooth coordinate function on  $\mathbb{R}^n$  to obtain

$$D_{v|_a}(x^j) \stackrel{\text{chain rule}}{=} v^i \frac{\partial}{\partial x^i}(x^j) \Big|_{x=a} = v^j = 0.$$

It follows that  $v_a$  is the zero vector so the map is injective.

To prove surjectivity, let  $w \in T_a \mathbb{R}^n$  be arbitrary. Let  $v = v^i e_i$ , where  $v^1, \dots, v^n \in \mathbb{R}$  are given by  $v^i = w(x^i)$ . We will show that  $w = D_{v|_a}$ . Let  $f$  be a smooth real-valued function on  $\mathbb{R}^n$ , by Taylor's theorem, we have

$$\begin{aligned} f(x) = & f(a) + \sum_{i=1}^n \frac{\partial f}{\partial x^i}(a)(x^i - a^i) \\ & + \sum_{i,j=1}^n (x^i - a^i)(x^j - a^j) \int_0^1 (1-t) \frac{\partial^2 f}{\partial x^i \partial x^j}(a + t(x-a)) dt. \end{aligned}$$

Then,

$$\begin{aligned} w(f) = & w(f(a)) + \sum_{i=1}^n w\left(\frac{\partial f}{\partial x^i}(a)(x^i - a^i)\right) \\ = & \sum_{i=1}^n \frac{\partial f}{\partial x^i}(a)(w(x^i) - w(a^i)) \\ = & \sum_{i=1}^n \frac{\partial f}{\partial x^i}(a)v^i = D_{v|_a} f. \end{aligned}$$

□

**Corollary 6.4.1.** For any  $a \in \mathbb{R}^n$ , the  $n$  derivations

$$\left. \frac{\partial}{\partial x^i} \right|_a, \dots, \left. \frac{\partial}{\partial x^n} \right|_a$$

defined by

$$\left. \frac{\partial}{\partial x^i} \right|_a (f) = \frac{\partial f}{\partial x^i}(a)$$

form a basis for  $T_a\mathbb{R}^n$ , which therefore has dimension  $n$ .

Now we can define tangent vectors on an arbitrary manifold. Let  $M$  be a manifold, and let  $p$  be a point of  $M$ . A linear map  $v : C^\infty(M) \rightarrow \mathbb{R}$  is called a *derivation* at  $p$  if it satisfies

$$v(fg) = f(p)v(g) + g(p)v(f) \quad \text{for all } f, g \in C^\infty(M).$$

**Definition 6.4.1.** (Tangent space to a manifold at a point) The set of all derivations of  $C^\infty(M)$  at  $p$ , denoted as  $T_pM$ , is a vector space called the *tangent space* to  $M$  at  $p$ . An element of  $T_pM$  is called a *tangent vector* at  $p$ .

The tangent vectors on manifolds have the same properties as that of Lemma 6.4.1. Recall that in Euclidean space, the total derivative of a smooth map at a point is a linear map that represents the best linear approximation to the map near the given point. In the case of manifolds, we will talk about the linear map between tangent spaces induced by smooth maps.

**Definition 6.4.2.** (Differential) Let  $M, N$  be manifolds and  $F : M \rightarrow N$  a smooth map, for each  $p \in M$  define a map

$$dF_p : T_pM \rightarrow T_{F(p)}N,$$

called the *differential* of  $F$  at  $p$ . Given  $v \in T_pM$ , we let  $dF_p(v)$  be the derivation at  $F(p)$  that acts on  $f \in C^\infty(N)$  by the rule

$$dF_p(v)(f) = v(f \circ F).$$

Note that  $f \circ F \in C^\infty(M)$  so  $v(f \circ F)$  is defined. The operator  $dF_p(v) : C^\infty(N) \rightarrow \mathbb{R}$  is linear because  $v$  is, and is a derivation at  $F(p)$  because for any  $f, g \in C^\infty(N)$ , we have

$$\begin{aligned} dF_p(v)(fg) &= v((fg) \circ F) = v((f \circ F)(g \circ F)) \\ &= (f \circ F)(p)v(g \circ F) + (g \circ F)(p)v(f \circ F) \\ &= (f \circ F)(p)dF_p(v)(g) + (g \circ F)(p)dF_p(v)(f). \end{aligned}$$

**Proposition 6.4.2.** Let  $M, N$ , and  $P$  be manifolds, let  $F : M \rightarrow N$  and  $G : N \rightarrow P$  be smooth maps, and let  $p \in M$ .

- (a)  $dF_p : T_pM \rightarrow T_{F(p)}N$  is linear.
- (b)  $d(G \circ F)_p = dG_{F(p)} \circ dF_p : T_pM \rightarrow T_{G \circ F(p)}P$ .
- (c)  $d(\text{Id}_M)_p = \text{Id}_{T_pM} : T_pM \rightarrow T_pM$ .
- (d) If  $F$  is a diffeomorphism, then  $dF_p : T_pM \rightarrow T_{F(p)}N$  is an isomorphism, and  $(dF_p)^{-1} = d(F^{-1})_{F(p)}$ .

Using corollary 6.4.1, we have preimages of the derivations (basis)  $\partial/\partial x^1|_{\varphi(p)}, \dots, \partial/\partial x^n|_{\varphi(p)}$  of  $T_{\varphi(p)}\mathbb{R}^n$  under the isomorphism  $d\varphi_p : T_pM \rightarrow T_{\varphi(p)}\mathbb{R}^n$  form a basis for  $T_pM$ .

Let  $F : M \rightarrow N$  be a smooth map between manifolds,  $(U, \varphi)$  a smooth chart of  $M$  containing  $p$  and  $(V, \psi)$  a smooth chart of  $N$  containing  $F(p)$ , we have  $\hat{F} = \psi \circ F \circ \varphi^{-1}$ ,  $\hat{p} = \varphi(p)$  then

$$dF_p \left( \left. \frac{\partial}{\partial x^i} \right|_p \right) = dF_p \left( d(\varphi^{-1})_{\hat{p}} \left( \left. \frac{\partial}{\partial x^i} \right|_{\hat{p}} \right) \right) = \frac{\partial \hat{F}^j}{\partial x^i}(\hat{p}) \left. \frac{\partial}{\partial y^j} \right|_{F(p)}$$

where  $(x^i)$ ,  $(y^j)$  is the local coordinate representation of the domain, codomain of  $F$ . The matrix of  $dF_p$  in terms of the coordinate basis or the Jacobian matrix of  $F$  is

$$\begin{pmatrix} \frac{\partial F^1}{\partial x^1}(p) & \dots & \frac{\partial F^1}{\partial x^n}(p) \\ \vdots & \ddots & \vdots \\ \frac{\partial F^m}{\partial x^1}(p) & \dots & \frac{\partial F^m}{\partial x^n}(p) \end{pmatrix}.$$

**Definition 6.4.3.** (Tangent bundle) Given a manifold  $M$ , we define the tangent bundle of  $M$ , denoted by  $TM$  by the disjoint union of the tangent spaces at all  $p \in M$ :

$$TM = \bigsqcup_{p \in M} T_pM$$

For any manifold  $M$ , the tangent bundle  $TM$  has a natural topology and smooth structure that makes it into a  $2n$ -dimensional smooth manifold and the projection  $\pi : TM \rightarrow M$  is smooth.

## 6.5 Submersions and Embeddings

In this section we will briefly talk about what does it mean for a manifold to be embedded in another manifold since for the purpose of the manifold assumption, we assume that the low dimensional manifold is embedded into a high dimensional Euclidean space.

**Definition 6.5.1.** (Submersion) A smooth map  $F : M \rightarrow N$  is called a smooth submersion if its differential is surjective at each point ( $\text{rank } F = \dim N$ ). It is called a smooth submersion if its differential is injective at each point ( $\text{rank } F = \dim M$ ).

**Definition 6.5.2.** (Embedding) A smooth immersion  $F : M \rightarrow N$  that is also a homeomorphism onto  $F(M) \subseteq N$  in the subspace topology is called a smooth embedding of  $M$  into  $N$ .

## 6.6 Covector fields and the differential of a function

**Definition 6.6.1.** (Covector) A linear functional  $\omega : V \rightarrow \mathbb{R}$  with  $V$  a finite-dimensional vector space is called a covector. The space of all covectors on  $V$ , called the dual space of  $V$ , is denoted as  $V^*$  and is itself a real vector space.

For a finite-dimensional vector space  $V$  with its basis  $(E_1, \dots, E_n)$ , let  $\varepsilon^1, \dots, \varepsilon^n \in V^*$  be the covectors defined by

$$\varepsilon^i(E_j) = \delta_j^i$$

where  $\delta_j^i$  is the Kronecker delta. Then it is easy to see that  $(\varepsilon^1, \dots, \varepsilon^n)$  forms a basis for  $V^*$ , called the dual basis to  $(E_1, \dots, E_n)$ . For example, denote  $(e^1, \dots, e^n)$  as the dual basis to the standard basis  $(e_1, \dots, e_n)$  in  $\mathbb{R}^n$  then for  $v \in \mathbb{R}^n$ , they can be given as

$$e^i(v) = e^i(v^1, \dots, v^n) = v^1 e^i(e_1) + \dots + v^n e^i(e_n) = v^i.$$

Thus, in general, if  $(E_j)$  is a basis for  $V$  and  $(\varepsilon^i)$  its dual basis then for any  $v = \sum_j v^j E_j := v^j E_j$ , we have

$$\varepsilon^i(v) := \varepsilon^i(v^j E_j) = \varepsilon^i \left( \sum_j v^j E_j \right) = \sum_j v^j \varepsilon^i(E_j) = \sum_j v^j \delta_j^i = v^i.$$

Furthermore, an arbitrary covector  $\omega \in V^*$  can be expressed in terms of the dual basis as

$$\omega = \omega(E_i) \varepsilon^i = \omega_i \varepsilon^i$$

and

$$\omega(v) = \omega(v^j E_j) = \omega_i \varepsilon^i(v^j E_j) = \omega_i v^i.$$

Let  $M$  be a manifold. For each  $p \in M$ , define the *cotangent space at  $p$*  to the tangent space  $T_p M$  to be  $T_p^* M = (T_p M)^*$ . Elements of  $T_p^* M$  are called *tangent covectors at  $p$* . Similar to the tangent bundle, denote the *cotangent bundle of  $M$*  as

$$T^* M = \bigsqcup_{p \in M} T_p^* M.$$

It has a natural projection map  $\pi : T^* M \rightarrow M$  sending  $\omega \in T_p^* M$  to  $p \in M$ . Given any smooth local coordinates  $(x^i)$  on  $U \subseteq M$ , for each  $p \in U$  we denote the basis for  $T_p^* M$  dual to  $(\partial/\partial x^i|_p)$  (the basis of  $T_p M$ ) by  $(\lambda^i|_p)$ . This defines  $n$  maps  $\lambda^1, \dots, \lambda^n : U \rightarrow T^* M$ , called *coordinate covector fields*. Given  $p \in U$ ,  $\lambda^i|_p \in T_p^* M$  picks out the  $i$ -th component of a tangent vector at  $p$ .

**Definition 6.6.2.** (Section) If  $\pi : M \rightarrow N$  is any continuous map, a section of  $\pi$  is a continuous right inverse for  $\pi$ , i.e., a continuous map  $\sigma : N \rightarrow M$  such that  $\pi \circ \sigma = \text{Id}_N$ .

A section of  $T^* M$  is called a *covector field*. Similarly, a section of the map  $\pi : TM \rightarrow M$  is called a *vector field on  $M$* . A vector field is called *smooth* if it is smooth as map from  $M$  to  $TM$  and the same definition applies for covector fields but from  $M$  to  $T^* M$ . A (co)vector field that is not necessarily smooth is called a rough (co)vector field.

In any smooth local coordinates on an open subset  $U \subseteq M$ , a (rough) covector field  $\omega$  can be written in terms of the coordinate covector fields  $(\lambda^i)$  as  $\omega = \omega_i \lambda^i$  for  $n$  functions  $\omega_i : U \rightarrow \mathbb{R}$  called the *component functions of  $\omega$*  characterized by

$$\omega_i(p) = \omega_p \left( \frac{\partial}{\partial x^i} \Big|_p \right) = \omega(p) \left( \frac{\partial}{\partial x^i} \Big|_p \right)$$

where  $w_p \in T_p^*M$ .

If  $\omega$  is a (rough) covector field and  $X$  is a vector field on  $M$ , then we can form a function  $\omega(X) : M \rightarrow \mathbb{R}$  by

$$\omega(X)(p) = \omega_p(X_p)$$

for  $p \in M$ .

The most important application of covector fields is providing a way to interpret partial derivatives as their components.

Let  $f$  be a smooth real-valued function on a manifold  $M$ . Define a covector field  $df$ , called the *differential of  $f$*  by

$$df_p(v) = v(f) \quad \text{for } v \in T_pM.$$

Let  $(x^i)$  be smooth coordinates on an open subset  $U \subseteq M$ , and let  $(\lambda^i|_p)$  be the basis for  $T_p^*M$ . Write  $df$  in coordinates as

$$df_p = A_i(p)\lambda^i|_p$$

for some component functions  $A_i : U \rightarrow \mathbb{R}$ , then we have

$$df_p \left( \frac{\partial}{\partial x^i} \Big|_p \right) = A_i(p)\lambda^i|_p \left( \frac{\partial}{\partial x^i} \Big|_p \right) = A_i(p).$$

By definition of  $df$ , it follows that

$$A_i(p) = df_p \left( \frac{\partial}{\partial x^i} \Big|_p \right) = \frac{\partial f}{\partial x^i}(p)$$

and

$$df_p = \frac{\partial f}{\partial x^i}(p)\lambda^i|_p.$$

Thus, the component functions of  $df$  in any smooth coordinate chart are partial derivatives of  $f$  w.r.t those coordinates. Because of this, we can think of  $df$  as an analogue of the classical gradient in  $\mathbb{R}^n$  in a coordinate-independent way. Let  $f = x^j$  be one of the coordinate functions then we obtain

$$dx^j|_p = \frac{\partial x^j}{\partial x^i}(p)\lambda^i|_p = \delta_i^j\lambda^i|_p = \lambda^j|_p.$$

Hence, the coordinate covector field is the differential of the coordinate functions,  $dx^j$ , and

$$df = \frac{df}{dx}dx.$$

It can be shown that  $df$  satisfies all of the usual properties of the ordinary derivative. Furthermore, let

$$\Delta f = f(p+v) - f(p)$$

for  $v \in \mathbb{R}^n$  then by Taylor's theorem,

$$\Delta f \approx \frac{\partial f}{\partial x^i}(p)v^i = df_p(v).$$

This shows that  $df_p$  is a linear functional that best approximates  $\Delta f$  near  $p$ .

In section 6.4, we defined  $df_p$  as a linear map from  $T_pM$  to  $T_{f(p)}\mathbb{R}$  but in this section, we defined it as a covector on  $T_pM$ , a linear map from  $T_pM$  to  $\mathbb{R}$  but they are the same because one can easily identify  $T_p\mathbb{R}$  with  $\mathbb{R}$ .

## 6.7 Tensors

**Definition 6.7.1.** (Multilinearity) Suppose  $V_1, \dots, V_k$  and  $W$  are vector spaces. A map  $F : V_1 \times \dots \times V_k \rightarrow W$  is said to be multilinear if it is linear as a function of each variable separately when the others are fixed.

Define  $L(V_1, \dots, V_k; W)$  as the set of all multilinear maps from  $V_1 \times \dots \times V_k$  to  $W$ , this is a vector space under pointwise addition and scalar multiplication. Let  $F_1, \dots, F_l \in L(V_1, \dots, V_k; W)$  depending on  $n_1, \dots, n_l$  variables then their *tensor product*  $F_1 \otimes \dots \otimes F_l$  is a multilinear function of  $n = n_1 + \dots + n_l$  variables such that

$$F_1 \otimes \dots \otimes F_l(v_1^{n_1}, \dots, v_{n_1}^{n_1}; v_1^{n_2}, \dots, v_{n_2}^{n_2}; \dots; v_1^{n_l}, \dots, v_{n_l}^{n_l}) = F_1(v_1^{n_1}, \dots, v_{n_1}^{n_1})F_2(v_1^{n_2}, \dots, v_{n_2}^{n_2}) \dots F_l(v_1^{n_l}, \dots, v_{n_l}^{n_l}).$$

Let  $V$  be a finite dimensional vector space. If  $k$  is a positive integer, a *covariant  $k$ -tensor on  $V$* ,  $\alpha$ , is an element of  $L(V1, \dots, V_k; \mathbb{R})$ , which can be thought of as a real-valued multilinear function of  $k$  elements of  $V$ :

$$\alpha : \underbrace{V \times \dots \times V}_{k \text{ times}} \rightarrow \mathbb{R}.$$

We will denote  $T^k(V^*)$  as  $L(V1, \dots, V_k; \mathbb{R})$ . A covariant  $k$ -tensor  $\alpha$  on  $V$  is said to be *symmetric* if its value is unchanged by interchanging any pair of arguments:

$$\alpha(v_1, \dots, v_i, \dots, v_j, \dots, v_k) = \alpha(v_1, \dots, v_j, \dots, v_i, \dots, v_k)$$

whenever  $1 \leq i < j \leq k$ . Denote the subspace of all symmetric  $k$ -tensors as  $\Sigma^k(V^*)$ . Let  $S_k$  be the symmetric group on  $k$  elements,  $\alpha$  a  $k$ -tensor, we then define a new  $k$ -tensor  $\sigma_\alpha$  by

$$\sigma_\alpha(v_1, \dots, v_k) = \alpha(v_{\sigma(1)}, \dots, v_{\sigma(k)}).$$

We define a projection  $\text{Sym} : T^k(V^*) \rightarrow \Sigma^k(V^*)$  called symmetrization by

$$\text{Sym } \alpha = (\text{Sym } \alpha)(v_1, \dots, v_k) = \frac{1}{k!} \sum_{\sigma \in S_k} \alpha(v_{\sigma(1)}, \dots, v_{\sigma(k)}) = \frac{1}{k!} \sum_{\sigma \in S_k} \sigma_\alpha.$$

If  $\alpha \in \Sigma^k(V^*)$  and  $\beta \in \Sigma^l(V^*)$ , we define their symmetric product to be the  $(k+l)$ -tensor  $\alpha\beta$  given by

$$\alpha\beta = \text{Sym } (\alpha \otimes \beta) = \frac{1}{(k+l)!} \sum_{\sigma \in S_{k+l}} \alpha(v_{\sigma(1)}, \dots, v_{\sigma(k)}) \beta(v_{\sigma(k+1)}, \dots, v_{\sigma(k+l)}).$$

**Proposition 6.7.1.** (a) *The symmetric product is symmetric and bilinear: for all symmetric tensors  $\alpha, \beta, \gamma$  and all  $a, b \in \mathbb{R}$ ,*

$$\alpha\beta = \beta\alpha,$$

$$(a\alpha + b\beta)\gamma = a\alpha\gamma + b\beta\gamma = \gamma(a\alpha + b\beta).$$

(b) *If  $\alpha$  and  $\beta$  are covectors, then*

$$\alpha\beta = \frac{1}{2}(\alpha \otimes \beta + \beta \otimes \alpha).$$

Alternatively, a covariant  $k$ -tensor  $\alpha$  on  $V$  is said to be *alternating* if it changes sign whenever two of its arguments are interchanged:

$$\alpha(v_1, \dots, v_i, \dots, v_j, \dots, v_k) = -\alpha(v_1, \dots, v_j, \dots, v_i, \dots, v_k).$$

The subspace of all alternating covariant  $k$ -tensors on  $V$  is denoted as  $\Lambda^k(V^*) \subset T^k(V^*)$ .

Let  $M$  be a manifold, we then define the bundle of covariant  $k$ -tensors on  $M$  by

$$T^k T^* M = \bigsqcup_{p \in M} T^k(T_p^* M)$$

and a tensor field is a smooth section of the canonical projection of the bundle. Let  $\alpha : M \rightarrow T^k T^* M$  be a  $k$ -tensor field then  $\alpha|_p \in T^k(T_p^* M) = L(\underbrace{T_p M, \dots, T_p M}_{k \text{ times}}; \mathbb{R})$ .

## 6.8 Riemannian Manifold

To be able to define geometric concepts such as lengths, angles, and distances on smooth manifolds, it is essential to introduce the structure of Riemannian metric.

**Definition 6.8.1.** (Riemannian manifold) Let  $M$  be a manifold, a Riemannian metric on  $M$  is a smooth symmetric covariant 2-tensor field on  $M$  that is positive definite at each point. A Riemannian manifold is a pair  $(M, g)$ , where  $M$  is a manifold and  $g$  a Riemannian metric on  $M$ .

If  $g$  is a Riemannian metric on  $M$ , then for each  $p \in M$ , the 2-tensor  $g_p$  is an inner product on  $T_p M$ . Because of this, we shall use the notation  $\langle v, w \rangle_g$  to denote  $g_p(v, w)$  for  $v, w \in T_p M$ .

In any smooth local coordinates  $(x^i)$  of an open subset  $U$  of  $M$ , a Riemannian metric  $g$  can be written

$$g = g_{ij} dx^i \otimes dx^j$$

where  $g_{ij} = g\left(\frac{\partial}{\partial x^i}, \frac{\partial}{\partial x^j}\right) : U \rightarrow \mathbb{R}$ . The symmetry of  $g$  allows us to write  $g$  as

$$\begin{aligned} g &= g_{ij} dx^i \otimes dx^j \\ &= \frac{1}{2}(g_{ij} dx^i \otimes dx^j + g_{ji} dx^i \otimes dx^j) \\ &= g_{ij} dx^i dx^j. \end{aligned}$$

**Example 6.8.1.** (Euclidean metric) Let  $g = \delta_{ij} dx^i dx^j$ ,  $p \in U \subseteq \mathbb{R}^n$  and  $v, w \in T_p M$  then

$$g_p(v, w) = \sum_{i=1}^n \sum_{j=1}^n \delta_{ij} dx^i|_p dx^j|_p(v, w) = \sum_{i=1}^n dx^i|_p(v) dx^i|_p(w) = \sum_{i=1}^n v^i w^i = v \cdot w$$

which is the Euclidean dot product.

The length or norm of  $v \in T_p M$  is defined to be

$$|v|_g = \langle v, v \rangle_g^{1/2}$$

and the angle between two nonzero tangent vectors  $v, w \in T_p M$  is a unique  $\theta \in [0, \pi]$  satisfying

$$\cos \theta = \frac{\langle v, w \rangle_g}{|v|_g |w|_g}.$$

One of the most important tools a Riemannian metric gives us is the ability to define lengths of curves. If  $M$  is a manifold, we define a curve in  $M$  to be a continuous map  $\gamma : J \rightarrow M$ , where  $J \subseteq \mathbb{R}$  is an interval. Let  $t_0 \in J$ , define the velocity of  $\gamma$  at  $t_0$ , denoted by  $\gamma'(t_0)$ , to be the vector

$$\gamma'(t_0) = d\gamma \left( \frac{d}{dt} \Big|_{t_0} \right) \in T_{\gamma(t_0)} M,$$

where  $d/dt|_{t_0}$  is the standard coordinate basis vector in  $T_{t_0} \mathbb{R}$ .

Let  $(U, \varphi)$  be a smooth chart with coordinate functions  $(x^i)$ . If  $\gamma(y_0) \in U$ , we can write the coordinate representation of  $\gamma$  as  $(\gamma^1(t), \dots, \gamma^n(t))$  then we have

$$\gamma'(t_0) = \frac{d\gamma^i}{dt}(t_0) \frac{\partial}{\partial x^i} \Big|_{\gamma(t_0)}$$

which is essentially the same formula as it would be in Euclidean space.

Given a curve  $\gamma : [a, b] \rightarrow M$ , the length of  $\gamma$  is

$$L_g(\gamma) = \int_a^b |\gamma'(t)|_g dt.$$

The integral is well-defined because  $|\gamma'(t)|_g$  is continuous for all but finitely many values of  $t$ , and has well defined limits from left, right end points.

**Proposition 6.8.1.** Let  $(M, g)$  be a Riemannian manifold, and let  $\gamma : [a, b] \rightarrow M$  be a piecewise smooth curve segment. If  $\tilde{\gamma}$  is a reparametrization of  $\gamma$ , then  $L_g(\tilde{\gamma}) = L_g(\gamma)$ .

Using curve segments as "measuring tapes", we can define distance between two points on a (connected) Riemannian manifold by letting it be the infimum of the lengths of all piece-wise smooth curve segments between two points.

## Manifold Regularization and Laplacian SVM

Regularization is a technique used to constraint a model to prevent overfitting. Manifold regularization is a somewhat similar process where you apply the regularization constraints with respect to the underlying manifold. We will first discuss the idea and theory behind some regularization techniques such as  $L_2$  regularization, commonly known as ridge regression, then connect it to manifold regularization and Laplacian SVM.

## 7.1 Regularization

Given an  $m \times n$  data matrix  $A$  with  $y$  a vector of labels associated with each row of  $A$ . In a usual Machine Learning setting, one would like to find a weight vector  $w$  such that  $Aw = y$  as a model to predict new data. But unfortunately, more often than not, the linear system  $Aw = y$  is ill-posed so we can not solve for an exact  $w$ . The obvious approach to this problem is to find  $w$  such that the (squared) error  $\|Aw - b\|_2^2$  is minimized.

Let  $f = w^T x$  be a linear hypothesis function for a certain linear regression problem. Note that if  $w$  is large then  $f$  would be sensitive to small perturbations in the input and thus, become undesirable as a classifier. To mitigate this problem, one would want to control the size of  $w$  by adding a regularization term, such as  $\|w\|_2^2$ , to the objective  $\|Aw - b\|_2^2$ . Our objective problem, often called ridge regression, is

$$\text{minimize} \quad \|Aw - y\|_2^2 + K \|w\|_2^2.$$

One may modify the objective problem by replacing the  $L_2$  regularization term with any other arbitrary  $L_p$  norm.

## 7.2 Manifold Regularization

Expanding the idea in the previous chapter to take into account the intrinsic geometry of the data, we arrive at manifold regularization which framework is given by [13]. Recall that in Machine Learning, the manifold assumption is an assumption that the data from an input space  $X$ , usually  $\mathbb{R}^n$ , only lies within a low dimensional manifold  $M \subset X$  which geometry will be used as a regularization term.

Let  $X = \{x_1, \dots, x_n\}$  be a dataset of  $n$  samples with each  $x_i$  drawn from a marginal distribution  $P_X$  with only  $l$  of which are labeled and drawn from the conditional distribution  $P(y | x)$ . The  $n - l = u$  unlabeled samples give additional information about the marginal distribution. Thus, there need to be an identifiable relation between the conditional and marginal for the model to work. To that end, we make a specific assumption that if two points  $x_1, x_2 \in X$  are close in the intrinsic geometry of  $P(X)$ , then the conditional distribution  $P(y | x_1)$  and  $P(y | x_2)$  are similar. We now proceed to a concrete formulation of the regularization problem.

## Results

## Conclusion

## Further Research

## Bibliography

- [1] Dirk Jan Struik. *Joseph-Louis Lagrange, comte de l'Empire — French mathematician — Britannica*. en. URL: <https://www.britannica.com/biography/Joseph-Louis-Lagrange-comte-de-lEmpire> (visited on 05/11/2022).
- [2] H. W. Kuhn and A. W. Tucker. “Nonlinear programming”. In: *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability, 1950*. Berkeley and Los Angeles: University of California Press, 1951, pp. 481–492.
- [3] Vladimir N. Vapnik. *An overview of statistical learning theory*. 1999. DOI: 10.1109/72.788640.
- [4] Vivek Salunkhe. *Support Vector Machine (SVM)*. 2021. URL: <https://medium.com/@viveksalunkhe80/support-vector-machine-svm-88f360ff5f38>.
- [5] Jocelyn Quaintance Jean Gallier. *Algebra, Topology, Differential Calculus, and Optimization Theory For Computer Science and Machine Learning*. 2022. URL: <https://www.cis.upenn.edu/~jean/math-deep.pdf>.
- [6] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag, 2006. ISBN: 0387310738.



- [7] Nello Cristianini and John Shawe-Taylor. *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*. 1st ed. Cambridge University Press, 2000. ISBN: 0521780195. URL: [http://www.amazon.com/Introduction-Support-Machines-Kernel-based-Learning/dp/0521780195/ref=sr\\_1\\_1?ie=UTF8&s=books&qid=1280243230&sr=8-1](http://www.amazon.com/Introduction-Support-Machines-Kernel-based-Learning/dp/0521780195/ref=sr_1_1?ie=UTF8&s=books&qid=1280243230&sr=8-1).
- [8] Chih-Chung Chang and Chih-Jen Lin Chih-Wei Hsu. “A Practical Guide to Support Vector Classification”. In: *BJU international* 101 (1 2008). ISSN: 1464-410X.
- [9] *Normed Vector Space of Rational Numbers is not Banach Space*. URL: [https://proofwiki.org/wiki/Normed\\_Vector\\_Space\\_of\\_Rational\\_Numbers\\_is\\_not\\_Banach\\_Space](https://proofwiki.org/wiki/Normed_Vector_Space_of_Rational_Numbers_is_not_Banach_Space).
- [10] Benyamin Ghogh et al. *Reproducing Kernel Hilbert Space, Mercer’s Theorem, Eigenfunctions, Nyström Method, and Use of Kernels in Machine Learning: Tutorial and Survey*. 2021. DOI: 10.48550/ARXIV.2106.08443. URL: <https://arxiv.org/abs/2106.08443>.
- [11] N. Aronszajn. “Theory of Reproducing Kernels”. In: *Transactions of the American Mathematical Society* 68 (3 1950). ISSN: 00029947. DOI: 10.2307/1990404.
- [12] J.M. Lee and J.M. Lee. *Introduction to Smooth Manifolds*. Graduate Texts in Mathematics. Springer, 2003. ISBN: 9780387954486. URL: <https://books.google.nl/books?id=eqfgZtjQceYC>.
- [13] Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. “Manifold Regularization: A Geometric Framework for Learning from Labeled and Unlabeled Examples”. In: *Journal of Machine Learning Research* 7.85 (2006), pp. 2399–2434. URL: <http://jmlr.org/papers/v7/belkin06a.html>.

# Appendix