

STOCKHOLM UNIVERSITY

BACHELOR THESIS

Useful Applications in Statistical Learning with Reproducing Kernel Hilbert Spaces

Author:
Yingjie CAO

Supervisor:
Yishao ZHOU

*A thesis submitted in fulfillment of the requirements
for the degree of Bachelor of Sciences*

in the

March 11, 2016

Declaration of Authorship

I, Yingjie CAO, declare that this thesis titled, “Useful Applications in Statistical Learning with Reproducing Kernel Hilbert Spaces” and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

STOCKHOLM UNIVERSITY

Abstract

Department of Mathematics

Bachelor of Sciences

Useful Applications in Statistical Learning with Reproducing Kernel Hilbert Spaces

by Yingjie CAO

This paper presents a general reproducing kernel Hilbert Spaces (RKHS) framework with its various applications in statistical learning area. This theory has been around for quite some time and has been widely used in nonlinear regression and classification problems. Kernel methods, which map data from low-dimensional space into higher-dimensional space (RKHS), can be transferred in many classical statistical learning algorithms. This paper can be roughly divided into two parts. In the first part, the writer attempts to take the reader from a very basic understanding of fields through Hilbert spaces, into reproducing kernel Hilbert spaces. In the second part, the writer want to show reader the abundant applications of kernel methods in statistical learning algorithms, with algorithms and read-world examples.

Keywords: RKHS, kernel methods, statistical learning, SVM

Acknowledgements

This project would not have been possible without the support of many people. Many thanks to my adviser, Yishao Zhou, who read my numerous revisions and helped make some sense of the confusion. Also thanks to my roommate, Jia Xu, who offered guidance and support. Thanks to School of Mathematical Sciences, Fudan University which offers me a great experience exchanging to Stockholm University. Thanks to the Department of Mathematics, Stockholm University which provides me comfortable and convenient research environment. And finally, thanks to my parents who endured this long process with me, always offering support and love.

Chapter 1

Introduction

1.1 Nonlinear Problem

Consider two concentric clouds of points (Figure 3.1), how can we distinguish them into two different groups with minimum error?

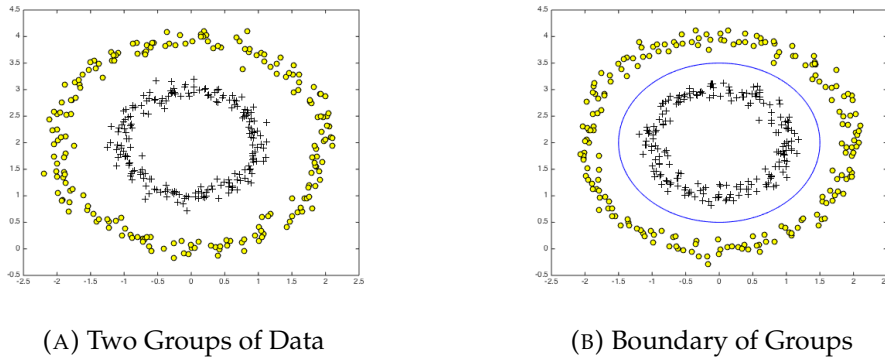


FIGURE 1.1: Example of Nonlinear Data

Of course, they are linear non-separable. Intuitively, we want to divide them with a circle, something like this (Figure 1.1b).

This case is simple, yet still a little tricky for the computer to solve it. The first thing is, we can no longer divide the two groups by lines. Then, we might think about adding features, from X_1, X_2 to $X_1, X_2, X_1^2, X_2^2, X_1X_2$. Therefore, the boundary can be written as this form:

$$a_1X_1 + a_2X_1^2 + a_3X_2 + a_4X_2^2 + a_5X_1X_2 + a_6 = 0 \quad (1.1)$$

The idea can be approached in another way. Let $Z_1 = X_1, Z_2 = X_1^2, Z_3 = X_2, Z_4 = X_2^2, Z_5 = X_1X_2$, (1.1) can be written as:

$$\sum_{i=1}^5 a_i Z_i + a_6 = 0 \quad (1.2)$$

Given the map $\varphi : \mathbb{R}^2 \rightarrow \mathbb{R}^5$, the two groups of data can be easily divided by the five dimensional hyperplane. In our simple example, we can simplify the five-dimension space into three dimensions, with $Z_1 = X_1^2, Z_2 = X_2^2, Z_3 = X_1X_2$, the two groups of data are shown as (Figure 1.2a). Through rotation (Figure 1.2b), the data can be easily separate by a plane.

Amazing? It is. Through the map $\varphi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$, linearly non-separable data groups can be separate by a hyperplane, and it is linear. Keep this point in mind and let's continue with another much more complex example.

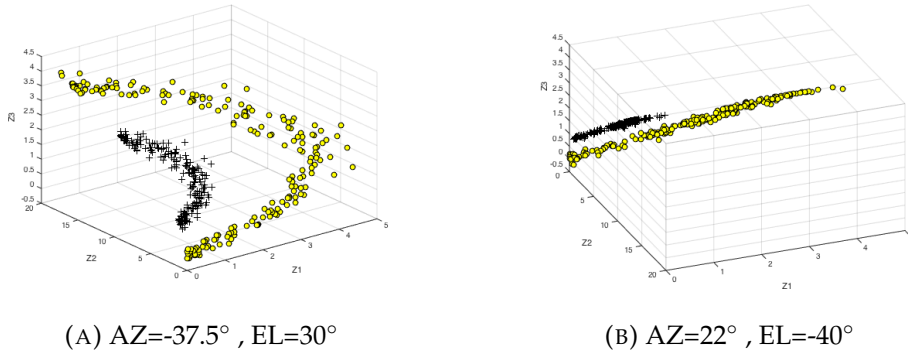


FIGURE 1.2: Nonlinear Data in 3D Space

1.2 Complex System

Adding features seems like a wise idea, yet it fails encountering with complex system. Another example of two groups of data are given below (Figure 1.3a), since the shape of boundary (Figure 1.3b) is no longer simple geometry, it is hard to tell how many polynomial features we should add. Though we can actually give an answer via method like Bias-Variance Trade-off, it still bothering and confusing.

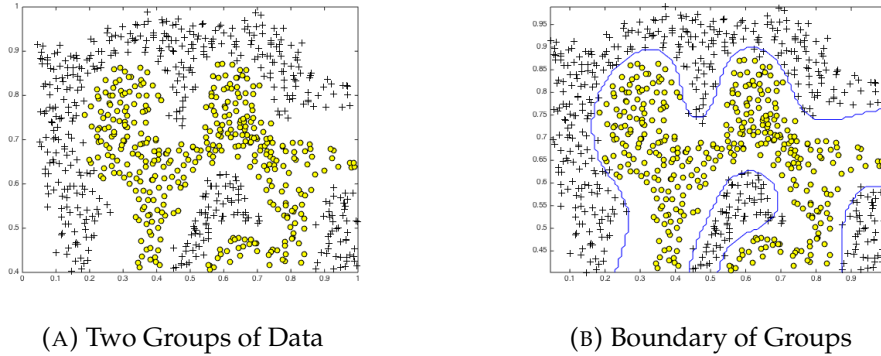


FIGURE 1.3: Example of Complex System

1.3 About this Paper

In Chapter 2, we introduce reproducing kernel Hilbert spaces (RKHS), which is the core idea in this paper. In Chapter 3, we use the Kernel Methods produced from RKHS in some statistical learning algorithms, and thus show the hidden power of this method. In Chapter 4, we introduce a real-world example, using Kernel Methods to solve it.

Chapter 2

Reproducing Kernel Hilbert Spaces

2.1 Regularization Problems

A general class of regularization problems has the form (see Appendix A)

$$\min_{f \in \mathcal{H}} \left[\sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f) \right] \quad (2.1)$$

where $L(y, f(x))$ is a loss function, $J(f)$ is a penalty functional, and \mathcal{H} is a space of functions on which $J(f)$ is defined.

2.2 Some Facts About RKHS

Let Ω be a set, for example, $\Omega = \mathbb{R}$, $\Omega = \{1, 2, \dots, N\}$, or $\Omega = [0, 1]$, or $\Omega = S_d$ (the d -dimensional sphere). A real symmetric function $K(x, y)$, $x, y \in \Omega$ is said to be positive definite on $\Omega \times \Omega$ if for every $n = 1, 2, \dots$, and every set of real numbers $\{a_1, a_2, \dots, a_N\}$ and $\{x_1, x_2, \dots, x_N\}$, $x_i \in \Omega$, we have $\sum_{i,j=1}^N a_i a_j K(x_i, x_j) \geq 0$. We have the famous

Theorem 2.2.1 (Moore-Aronszajn Theorem). *To every positive definite function K on $\Omega \times \Omega$ there corresponds a unique RKHS \mathcal{H}_K of real valued functions on Ω and vice versa.*

Given K , we can construct \mathcal{H}_K in the form of all finite linear combinations $\sum_{i=1}^N a_i K_{x_i}$ with $x_i \in \Omega$ and limits of such function as the x_i become dense in Ω , in the norm induced by the inner product

$$\langle K(\cdot, x_i), K(\cdot, x_j) \rangle_{\mathcal{H}_K} = K(x_i, x_j) \quad (2.2)$$

The function $K(\cdot, x_i)$ is so-called representer of evaluation at x_i in \mathcal{H}_K , that is to say, for any $f \in \mathcal{H}_K$ and fixed x_i ,

$$\langle K(\cdot, x_i), f \rangle_{\mathcal{H}_K} = f(x_i) \quad (2.3)$$

If $K(x, y)$ has a representation of the form

$$K(x, y) = \sum_i \gamma_i \phi_i(x) \phi_i(y) \quad (2.4)$$

with

$$\int_{\Omega} \phi_{\xi}(s) \phi_{\eta}(s) d\mu(s) = \begin{cases} 1, & \xi = \eta \\ 0, & \text{otherwise} \end{cases} \quad (2.5)$$

where μ is some measure on Ω , then

$$\langle f, g \rangle_{\mathcal{H}_K} = \sum_i \frac{f_i g_i}{\gamma_i} \quad (2.6)$$

where $f_i = \int \phi_i(s) f(s) d\mu(s)$ and similarly for g_i .

In particular,

$$\langle \phi_{\xi}, \phi_{\eta} \rangle_{\mathcal{H}_K} = \begin{cases} \frac{1}{\gamma_{\xi}}, & \xi = \eta \\ 0, & \text{otherwise} \end{cases} \quad (2.7)$$

$K(\cdot, \cdot)$ is known as the reproducing kernel (RK) for \mathcal{H}_K , due to the "reproducing property".

Corollary 2.2.1.1. *If $g(x) = \sum_{i=1}^N \alpha_i K(x, x_i)$, then*

$$\begin{aligned} \|g\|_{\mathcal{H}_K}^2 &= \langle g, g \rangle_{\mathcal{H}_K} = \langle \alpha_i K(\cdot, x_i), \alpha_j K(\cdot, x_j) \rangle_{\mathcal{H}_K} \\ &= \sum_{i=1}^N \sum_{j=1}^N \langle K(\cdot, x_i), K(\cdot, x_j) \rangle_{\mathcal{H}_K} \alpha_i \alpha_j \\ &= \sum_{i=1}^N \sum_{j=1}^N K(x_i, x_j) \alpha_i \alpha_j \end{aligned} \quad (2.8)$$

Theorem 2.2.2 (The Representer Theorem). *Let $\{\phi_i\}_{i=1}^M$ be M functions on Ω^2 with the property that the $M \times N$ matrix T with ij th entry $\phi_i(x_j)$ is of rank M . Let $L(y_j, f)$ be a functional of f which depends on f only through $f(x_j) = f_j$. Then the solution to the problem: find $f \in \text{span}\{\phi_i\} + h$ with $h \in \mathcal{H}_K$ to minimize*

$$\sum_{i=1}^N L(y_j, f(x_j)) + \lambda \|h\|_{\mathcal{H}_K}^2 \quad (2.9)$$

has a representation of the form

$$f(\cdot) = \sum_{i=1}^M c_i \phi_i(\cdot) + \sum_{j=1}^N d_j K(\cdot, x_j) \quad (2.10)$$

Corollary 2.2.2.1. *Suppose that $\tilde{g}(x) = g(x) + \rho(x)$, with $\rho(x) \in \mathcal{H}_K$, and orthogonal in \mathcal{H}_K to each of $K(x, x_i), i = 1, \dots, N$. Then*

$$\sum_{i=1}^N L(y_i, \tilde{g}(x_i)) + \lambda J(\tilde{g}) \geq \sum_{i=1}^N L(y_i, g(x_i)) + \lambda J(g) \quad (2.11)$$

with equation iff $\rho(x) = 0$.

2.3 Spaces of Functions Generated by Kernels

An important subclass of problems of the form (2.1) are generated by a positive definite kernel $K(x, y)$, and the corresponding space of functions \mathcal{H}_K as a reproducing kernel Hilbert space. The penalty functional J is defined in terms of the kernel as well.

Let $x, y \in \mathbb{R}^p$. We consider the space of functions generated by the linear span of $K(\cdot, y)$, $y \in \mathbb{R}^p$. Suppose that K has an eigen-expansion

$$K(x, y) = \sum_{i=1}^{\infty} \gamma_i \phi_i(x) \phi_i(y) \quad (2.12)$$

with $\gamma_i \geq 0$, $\sum_{i=1}^{\infty} \gamma_i^2 \leq \infty$. Elements of \mathcal{H}_K have an expansion in terms of these eigen-functions,

$$f(x) = \sum_i c_i \phi_i(x) \quad (2.13)$$

According to (2.6), the norm induced by K can be shown as

$$\|f\|_{\mathcal{H}_K}^2 = \langle f, f \rangle_{\mathcal{H}_K} = \sum_{i=1}^{\infty} \frac{f_i^2}{\gamma_i} = \sum_{i=1}^{\infty} \frac{c_i^2}{\gamma_i} < \infty \quad (2.14)$$

where

$$f_i = \int \phi_i(s) f(s) d\mu(s) = \int \phi_i(s) \left(\sum_j c_j \phi_j(s) \right) d\mu(s) = c_i \quad (2.15)$$

The penalty functional in (2.1) for the space \mathcal{H}_K is defined to be the squared norm $J(f) = \|f\|_{\mathcal{H}_K}^2$. The quantity $J(f)$ can be interpreted as a generalized ridge penalty, where functions with large eigenvalues in the expression (2.14) get penalized less, and vice versa.

Rewriting (2.1) we have

$$\min_{f \in \mathcal{H}_K} \left[\sum_{i=1}^N L(y_i, f(x_i)) + \lambda \|f\|_{\mathcal{H}_K}^2 \right] \quad (2.16)$$

or equivalently,

$$\min_{\{c_j\}_1^{\infty}} \left[\sum_{i=1}^N L(y_i, \sum_{j=1}^{\infty} c_j \phi_j(x_i)) + \lambda \sum_{j=1}^{\infty} \frac{c_j^2}{\gamma_j} \right] \quad (2.17)$$

Accordance to Theorem 2.2.2 and Corollary 2.2.2.1, the solution for (2.16) can be shown as finite-dimensional, and has the form

$$f(x) = \sum_{i=1}^N \alpha_i K(x, x_i). \quad (2.18)$$

with Corollary 2.2.1.1,

$$\|f\|_{\mathcal{H}_K}^2 = \sum_{i=1}^N \sum_{j=1}^N K(x_i, x_j) \alpha_i \alpha_j. \quad (2.19)$$

Therefore, the regularization problem can be simplified into a finite-dimensional criterion

$$\min_{\alpha} L(\mathbf{y}, \mathbf{K}\alpha) + \lambda \alpha^T \mathbf{K} \alpha \quad (2.20)$$

where α is a vector with i th entry $\alpha(i)$, \mathbf{K} is a $N \times N$ matrix with ij th entry $K(x_i, x_j)$.

2.4 Custom Kernel

When we are dealing with regularization problems with kernel methods, it is always important to choose appropriate kernel functions, since the performance of kernel algorithm is highly dependent on the kernel. Though choosing the best kernel function depends on the specific problem, it is better to have prior knowledge about some most commonly used kernels.

Below is a list of some kernel functions available from the existing literature.

1. Linear Kernel

The Linear kernel is the simplest kernel function. Kernel algorithms using a linear kernel are often equivalent to their non-kernel counterparts.

$$k(x, y) = x^T y + c$$

2. Polynomial Kernel

The Polynomial kernel is a non-stationary kernel. Polynomial kernels are well suited for problems where all the training data is normalized.

$$k(x, y) = (\alpha x^T y + c)^d$$

3. Radial Basis Function Kernel

(a) Gaussian Kernel

It is the most commonly used kernel function and also acts well in many cases.

$$k(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right)$$

The adjustable parameter σ plays a major role in the performance of the kernel, and should be carefully tuned to the problem at hand. If overestimated, the exponential will behave almost linearly and the higher-dimensional projection will start to lose its non-linear power. In the other hand, if underestimated, the function will lack regularization and the decision boundary will be highly sensitive to noise in training data.

(b) Exponential Kernel

The exponential kernel is closely related to the Gaussian kernel, with only the square of the norm left out.

$$k(x, y) = \exp\left(-\frac{\|x - y\|}{2\sigma^2}\right)$$

(c) Laplacian Kernel

The Laplace Kernel is completely equivalent to the exponential kernel, except for being less sensitive for changes in the sigma parameter.

$$k(x, y) = \exp\left(-\frac{\|x - y\|}{\sigma}\right)$$

(d) ANOVA Kernel

The ANOVA kernel is said to perform well in multidimensional regression problems (Hofmann, 2008).

$$k(x, y) = \sum_{k=1}^n \exp(-\sigma(x^k - y^k)^2)^d$$

4. Hyperbolic Tangent (Sigmoid) Kernel

The Sigmoid Kernel comes from the Neural Networks field, where the bipolar sigmoid function is often used as an activation function for artificial neurons. Since SVM model using a sigmoid kernel function is equivalent to a two-layer, perceptron neural network, this kernel was quite popular for support vector machines. Also, despite being only conditionally positive definite, it has been found to perform well in practice.

$$k(x, y) = \tanh(\alpha x^T + c)$$

There are two adjustable parameters in the sigmoid kernel, the slope α and the intercept constant c . A common value for α is $1/N$, where N is the data dimension.

5. Log Kernel

The Log kernel seems to be particularly interesting for images, but is only conditionally positive definite.

$$k(x, y) = -\log(\|x - y\|^d + 1)$$

6. Rational Quadratic Kernel

The Rational Quadratic kernel is less computationally intensive than the Gaussian kernel and can be used as an alternative when using the Gaussian becomes too expensive.

$$k(x, y) = 1 - \frac{\|x - y\|^2}{\|x - y\|^2 + c}$$

7. Cauchy Kernel

The Cauchy kernel comes from the Cauchy distribution (Basak, 2008). It is a long-tailed kernel and can be used to give long-range influence and sensitivity over the high dimension space.

$$k(x, y) = \frac{1}{1 + \frac{\|x - y\|^2}{\sigma^2}}$$

Chapter 3

Applications in Different Statistical Learning Algorithms

3.1 Kernel PCA

Standard linear principal components (PCA) are obtained from the eigenvectors of the covariance matrix, and give directions in which the data have maximal variance. Kernel PCA is the method to expand the features by non-linear transformations, and then apply PCA in this transformed feature space.

Lemma 3.1.1. *Given a data matrix X , the inner-product (gram) matrix $K = XX^T = \langle x_i, x_j \rangle$. The double-centered version of the gram matrix can be written as:*

$$\tilde{K} = (I - \mathbf{1}_N)K(I - \mathbf{1}_N) = UD^2U^T \quad (3.1)$$

with $\mathbf{1}_N$ a $N \times N$ matrix for which each element takes value $1/N$.

Proof. Let $K = XX^T = \langle x_i, x_j \rangle$,

$$\begin{aligned} \Delta_{ij}^2 &= \|x_i - x_j\|^2 = \|(x_i - \bar{x}) - (x_j - \bar{x})\|^2 \\ &= \|x_i - \bar{x}\|^2 + \|x_j - \bar{x}\|^2 - 2\langle x_i - \bar{x}, x_j - \bar{x} \rangle \end{aligned} \quad (3.2)$$

Then we have double-centered matrix \tilde{K}

$$\begin{aligned} \tilde{K} &= \langle x_i - \bar{x}, x_j - \bar{x} \rangle \\ &= \frac{1}{2} \left[\|x_i - \bar{x}\|^2 + \|x_j - \bar{x}\|^2 - \Delta_{ij}^2 \right] \\ &= \langle x_i, x_j \rangle - \frac{2}{N} \langle x_i, x_j \rangle + \frac{1}{N^2} \langle x_i, x_j \rangle \\ &= K - \mathbf{1}_N K - K \mathbf{1}_N + \mathbf{1}_N K \mathbf{1}_N \\ &= (I - \mathbf{1}_N)K(I - \mathbf{1}_N) \end{aligned} \quad (3.3)$$

The centered version of X can be decomposed by its SVD in the form $\tilde{X} = (I - \mathbf{1}_N)X = UDV^T$, then $Z = UD$ is the matrix of principal components variables. And it follows that $\tilde{K} = \tilde{X}\tilde{X}^T = UDV^TVD^T U^T = UD^2U^T$, we can compute Z from the eigen decomposition of K . \square

Kernel PCA simply mimics this procedure, interpreting the kernel matrix $K = \{K(x_i, x_j)\}$ as an inner-product matrix of the implicit features $\langle \phi(x_i), \phi(x_j) \rangle$ and finding its eigenvectors. The elements of the m th component z_m (m th column of Z) can be written as $z_{im} = \sum_{j=1}^N \alpha_{jm} K(x_i, x_j)$, where $\alpha_{jm} = u_{jm}/d_m$. Now we will show the process of derivation.

Assume that our data mapped into feature space, $\phi(x_1), \dots, \phi(x_N)$, is centered, i.e. $\sum_{k=1}^N \phi(x_k) = 0$ (if not, using the formula of centralization by lemma 3.1.1). To do PCA for the covariance matrix

$$C = \sum_{i=1}^N \phi(x_i) \phi(x_i)^T \quad (3.4)$$

we have to find eigenvalues $\lambda \geq 0$ and eigenvectors $V \in \mathcal{F} \setminus \{0\}$ satisfying $\lambda V = CV$. Note that all solutions V lie in the span of $\phi(x_1), \dots, \phi(x_N)$, we may consider the equivalent system

$$\lambda(\phi(x_k) \cdot V) = (\phi(x_k) \cdot CV) \text{ for all } k = 1, \dots, N. \quad (3.5)$$

and there exist coefficients $\alpha_1, \dots, \alpha_N$ such that

$$V = \sum_{i=1}^N \alpha_i \phi(x_i) \quad (3.6)$$

considering the $N \times N$ matrix K with ij th entry

$$K_{ij} = \langle \phi(x_i), \phi(x_j) \rangle = (\phi(x_i) \cdot \phi(x_j)) \quad (3.7)$$

we arrive at

$$N\lambda K\alpha = K^2\alpha \quad (3.8)$$

where α denotes the column vector with entries $\alpha_1, \dots, \alpha_N$. To find solutions of (3.8), we solve the eigenvalue problem

$$N\lambda\alpha = K\alpha \quad (3.9)$$

for nonzero eigenvalues. Clearly, all solutions of (3.8) do satisfy (3.9).

We normalize the solutions α^k belonging to nonzero eigenvalues by requiring that the corresponding vectors in F be normalized, i.e. $(V^k \cdot V^k) = 1$. By virtue of (3.6), (3.7) and (3.9), this translates into

$$1 = \sum_{i,j=1}^N \alpha_i^k \alpha_j^k (\phi(x_i) \cdot \phi(x_j)) = (\alpha^k \cdot K\alpha^k) = \lambda_k (\alpha^k \cdot \alpha^k) \quad (3.10)$$

For principle component extraction, we compute projections of the image of a test point $\phi(x)$ onto the eigenvectors V^k in \mathcal{F} accordance to

$$(V^k \cdot \phi(x)) = \sum_{i=1}^N \alpha_i^k (\phi(x_i) \cdot \phi(x)) \quad (3.11)$$

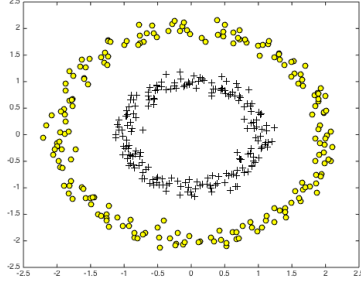
Consider two concentric clouds of points (has been shown in Introduction). We wish to use kernel PCA to identify these groups. Here we use two kinds of kernel:

- Polynomial Kernel:

$$k(x, y) = (\alpha x^T y + c)^d$$

- Gaussian Kernel:

$$k(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right)$$



(A) Input points before kernel PCA

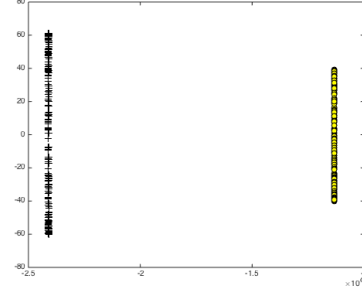
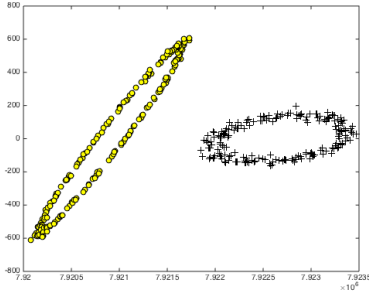
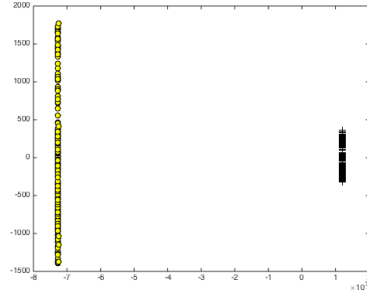
(B) Output points after kernel PCA with Gaussian Kernel $k(x, y) = \exp(-\|x - y\|^2)$ (C) Output points after kernel PCA with Linear Kernel $k(x, y) = x^T y + 1$ (D) Output points after kernel PCA with Polynomial Kernel $k(x, y) = (x^T y + 1)^2$

FIGURE 3.1: Kernel PCA

3.2 Kernel SVM

3.2.1 Support Vector Classifier

The left panel (3.2a) shows the separable case. The linear decision boundary is $x^T \beta + \beta_0 = 0$, which bound the maximal margin of width $2M = 2/\|\beta\|$, where the margin been determined by the points with its margin width M for one side. The right panel (3.2b) shows the non-separable (overlap) case. Some points are on the wrong side of their margin by ξ_i ; points on the correct side have $\xi_i = 0$. The margin is maximized subject to a total budget $\sum \xi_i \leq \text{constant}$.

Our training data consists of N pairs $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$, with $x_i \in \mathbb{R}^p$ and $y_i \in \{-1, 1\}$. Define a hyperplane by

$$\{x : f(x) = x^T \beta + \beta_0 = 0\} \quad (3.12)$$

where β is a unit vector: $\|\beta\| = 1$. A classification rule induce by $f(x)$ is

$$G(x) = \text{sign}[x^T \beta + \beta_0] \quad (3.13)$$

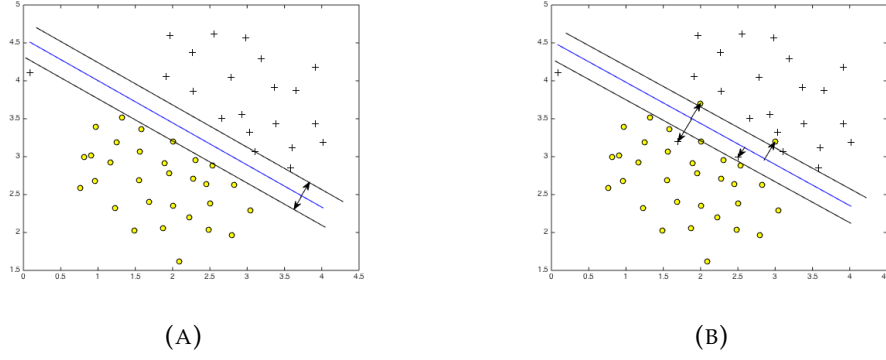


FIGURE 3.2: Support Vector Classifier

For the separable case (3.2a), we can find a function $f(x) = x^T \beta + \beta_0$ with $y_i f(x_i) > 0 \forall i$. The optimization problem can be shown as

$$\begin{aligned} & \max_{\beta, \beta_0, \|\beta\|=1} M \\ & \text{subject to } y_i(x_i^T \beta + \beta_0) \geq M, i = 1, \dots, N \end{aligned} \quad (3.14)$$

or equivalently,

$$\begin{aligned} & \min_{\beta, \beta_0} \|\beta\| \\ & \text{subject to } y_i(x_i^T \beta + \beta_0) \geq 1, i = 1, \dots, N \end{aligned} \quad (3.15)$$

Suppose now that the classes overlap in feature space (3.2b), we can modify (3.15) with a little change

$$\begin{aligned} & \min_{\beta, \beta_0} \|\beta\| \\ & \text{subject to } \begin{cases} y_i(x_i^T \beta + \beta_0) \geq 1 - \xi_i, \forall i \\ \xi_i \geq 0, \sum \xi_i \leq \text{constant} \end{cases} \end{aligned} \quad (3.16)$$

or equivalently,

$$\begin{aligned} & \min_{\beta, \beta_0} \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i \\ & \text{subject to } \xi_i \geq 0, y_i(x_i^T \beta + \beta_0) \geq 1 - \xi_i, \forall i \end{aligned} \quad (3.17)$$

The Lagrange function is

$$L_P = \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [y_i(x_i^T \beta + \beta_0) - (1 - \xi_i)] - \sum_{i=1}^N \mu_i \xi_i \quad (3.18)$$

which we minimize w.r.t β , β_0 and ξ_i . Setting the respective derivatives to zero, we get

$$\beta = \sum_{i=1}^N \alpha_i y_i x_i, \quad (3.19)$$

$$0 = \sum_{i=1}^N \alpha_i y_i, \quad (3.20)$$

$$\alpha_i = C - \mu_i, \forall i \quad (3.21)$$

as well as the positivity constraints $\alpha_i, \mu_i, \xi_i \geq 0 \forall i$. By substituting (3.19)-(3.20) into (3.18), we obtain the Lagrangian dual object function

$$L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j \quad (3.22)$$

$$\text{subject to } 0 \leq \alpha_i \leq C, \text{ and } \sum_{i=1}^N \alpha_i y_i = 0$$

From (3.19) we see that the solution for β has the form

$$\hat{\beta} = \sum_{i=1}^N \hat{\alpha}_i y_i x_i \quad (3.23)$$

Given the solutions $\hat{\beta}_0$ and $\hat{\beta}$, the decision function can be written as

$$G(\hat{x}) = \text{sign}[\hat{f}(x)] = \text{sign}[x^T \hat{\beta} + \hat{\beta}_0] \quad (3.24)$$

3.2.2 Support Vector Machines and Kernels

The support vector classifier described so far finds linear boundaries in the input feature space. Generally linear boundaries in the enlarged space achieve better training-class separation, and translate to nonlinear boundaries in the original space. Once the basis functions $h_m(x)$, $m = 1, \dots, M$ are selected, the procedure is the same as before. We fit the SV classifier using input features $h(x_i) = (h_1(x_i), h_2(x_i), \dots, h_M(x_i))$, $i = 1, \dots, N$, and produce the nonlinear function $\hat{f}(x) = h(x)^T \hat{\beta} + \hat{\beta}_0$. The classifier is $\hat{G}(x) = \text{sign}(\hat{f}(x))$ as before.

The support vector machine (SVM) classifier is an extension of this idea, where the dimension of the enlarged space is allowed to get very large, infinite in some cases. It might seem that the computations would become prohibitive. It would also seem that with sufficient basis functions, the data would be separable, and overfitting would occur.

We can represent the optimization problem (3.18) and its solution in a special way that only involves the input features via inner products. We do this directly for the transformed feature vectors $h(x_i)$. We then see that for particular choices of h , these inner products can be computed very cheaply.

With prior knowledge of the kernel function

$$K(x, y) = \langle h(x), h(y) \rangle \quad (3.25)$$

The Lagrange dual function (3.22) has the form

$$\begin{aligned} L_D &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle h(x_i), h(x_j) \rangle \\ &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j K(x_i, x_j) \end{aligned} \quad (3.26)$$

The solution function $f(x)$ can be written as

$$\begin{aligned} f(x) &= h(x)^T \beta + \beta_0 \\ &= \sum_{i=1}^N \alpha_i y_i \langle h(x), h(x_i) \rangle + \beta_0 \\ &= \sum_{i=1}^N \alpha_i y_i K(x, x_i) + \beta_0 \end{aligned} \quad (3.27)$$

As before, given α_i, β_0 can be determined by solving $y_i f(x_i) = 1$ in (3.26) for all x_i for which $0 < \alpha_i < C$.

With $f(x) = h(x)^T \beta + \beta_0$, we have

$$\begin{cases} \xi_i \geq 0 \\ y_i f(x_i) \geq 1 - \xi_i \end{cases} \quad \forall i \quad (3.28)$$

Therefore, the optimization problem (3.17) can be rewritten as the form

$$\min_{\beta, \beta_0} \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N [1 - y_i f(x_i)]_+ \quad (3.29)$$

where the subscript "+" indicates positive part. With $\lambda = 1/C$, we have

$$\min_{\beta, \beta_0} \sum_{i=1}^N [1 - y_i f(x_i)]_+ + \frac{\lambda}{2} \|\beta\|^2 \quad (3.30)$$

This has the form with a loss function and a shrinkage, which is a familiar paradigm in function estimation. Next, we will use the knowledge of reproducing kernel Hilbert spaces to solve this regularization problem.

3.2.3 Reproducing Kernels in SVM

Suppose the basis h arises from the eigen-expansion of a positive definite kernel K ,

$$K(x, y) = \sum_{j=1}^{\infty} \gamma_j \phi_j(x) \phi_j(y) \quad (3.31)$$

and $h_j(x) = \sqrt{\gamma_j} \phi_j(x)$. Then with $\theta_j = \sqrt{\gamma_j} \beta_j$, we can write (3.30) as

$$\min_{\beta_0, \{\theta_j\}_1^{\infty}} \sum_{i=1}^N \left[1 - y_i (\beta_0 + \sum_{j=1}^{\infty} \theta_j \phi_j(x_i)) \right]_+ + \frac{\lambda}{2} \sum_{j=1}^{\infty} \frac{\theta_j^2}{\gamma_j} \quad (3.32)$$

The theory of reproducing kernel Hilbert spaces described there guarantees a finite-dimensional solution of the form

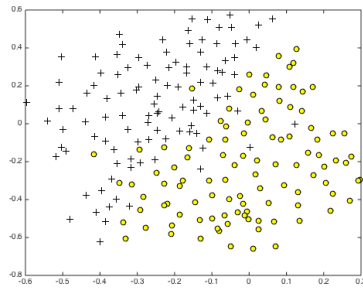
$$f(x) = \beta_0 + \sum_{i=1}^N \alpha_i K(x, x_i) \quad (3.33)$$

Also the equivalent version of the optimization criterion of (3.32) is

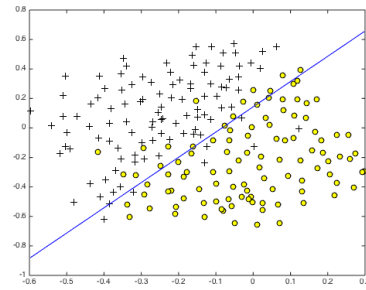
$$\min_{\beta_0, \theta} \sum_{i=1}^N (1 - y_i f(x_i))_+ + \frac{\lambda}{2} \alpha^T K \alpha \quad (3.34)$$

where K is the $N \times N$ matrix of kernel evaluations for all pairs of training features.

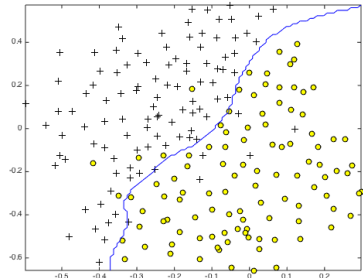
Here is an example that classes overlap in feature space, we will see how the decision boundary changes with different kind of kernels.



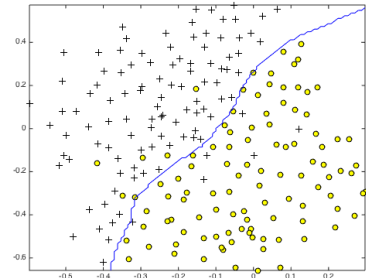
(A) Input points: non-separable case



(B) Training SVM with Linear Kernel
 $k(x, y) = x^T y + 1$



(C) Training SVM with Gaussian
Kernel $k(x, y) = \exp(-\frac{\|x-y\|^2}{2*0.1^2})$



(D) Training SVM with Log Kernel
 $k(x, y) = -\log(\|x-y\|^{0.5} + 1)$

FIGURE 3.3: Training SVM with Different Kernels

Chapter 4

Real-World Example: handwritten digits recognition¹

The MNIST database of handwritten digits has a training set of 60,000 examples, and a test set of 10,000 examples. The digits have been size-normalized and centered in a fixed-size image, so it is easy to apply our supervised statistical learning algorithms with minimal efforts on preprocessing and formatting. We want to compare the effects among different algorithms, with and without kernel. Let's have a overview towards these digits, here (4.1) are the first 100 digits in training dataset.

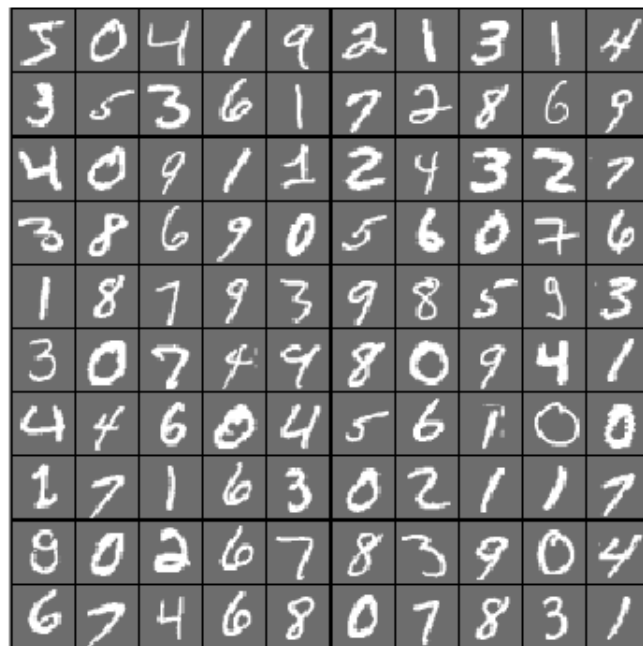


FIGURE 4.1: MNIST database: handwritten digits

4.1 SVM: No kernel v.s. Kernel

LIBSVM is a powerful software for support vector classification, it uses one v.s. one strategy to deal with multi-class classification problem. We use

¹All the code and data used in this chapter can be find here: <https://github.com/Louise222/Bachelor-Thesis.git>

this package to train two models, one without kernel and one with gaussian kernel. By comparing the predicted label and real label in the testing dataset, we can get accuracy of these two models.

Before training model, we need to preprocess the data. Our training dataset is large, it is time consuming either training model or searching for a better parameters. Thus, we subset the training datasets into two parts, and only use 1000 digits among 60,000 to search for good parameters. Later in training model process, we still use the whole training dataset.

```
% step1: turn into LIBSVM format
% <label> <index1>:<value1> <index2>:<value2> ...
train_sparse = sparse(train_images_scale);
libsvmwrite('train',train_labels,train_sparse);
% step2: ../tool/subset.py for a subset of data
>subset.py train 1000 output1 output2
% step3: ../tool/grid.py for searching parameters
>grid.py -log2c -2,10,1 -log2g -4,-14,-1 output1
```

This figure (4.2) is plotted by **gnuplot** followed by ">grid.py" command. It finds the best parameters for SVM with RBF kernel.

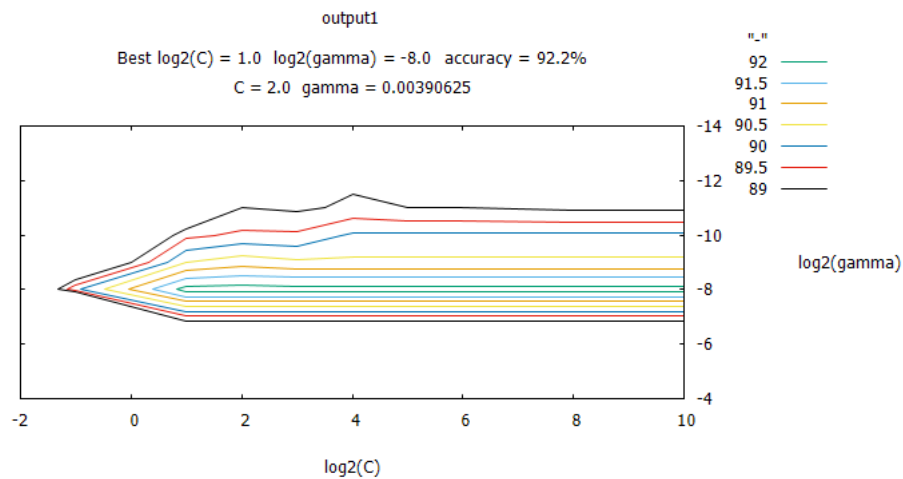


FIGURE 4.2: Searching for the Best Parameters

In our example, we get $C = 2$, $g = 0.0039$, that is to say, it is a good idea to set $C = 2$ for shrinkage parameter and $\gamma = 0.0039$ for Gaussian kernel parameter.

```
% Train and Predict, linear kernel
modell = svmtrain(train_labels, train_images_scale, ...
    '-s 0 -t 0 -c 2');
[predicted_label1, accuracy1, decision_values1] = ...
    svmpredict(test_labels, test_images_scale, modell);
% Shown wrongly-matched digits
index1 = find(test_labels ~= predicted_label1);
display_network(test_images(index1(1:100), :));
```

Without kernel, we get accuracy of 93.16%.

```
% Train and Predict, gaussian kernel
model2 = svmtrain(train_labels, train_images_scale,...
    '-s 0 -t 2 -g 0.0039 -c 2');
[predicted_label2, accuracy2, decision_values2] =...
    svmpredict(test_labels, test_images_scale, model2);
% Shown wrongly-matched digits
index2 = find(test_labels~=predicted_label2);
display_network(test_images(index2(1:100),:))';
```

With kernel, we get accuracy of 98.35%.

Here are some wrongly-matched digits predicted by SVM. The digits in the left panel (4.3a) are wrongly predicted by linear SVM, and those in the right panel (4.3b) are wrongly predicted by kernel SVM.



(A) No Kernel



(B) Gaussian Kernel

FIGURE 4.3: Wrongly-matched Digits by SVM

In conclusion, we have 5.84% of predicting error with linear SVM, and only have 1.65% of predicting error with kernel SVM.

4.2 Neural Network

Neural Network method uses sigmoid kernel. In fact, support vector classifier with a sigmoid kernel acts exactly with a two-layer neural network model. In this section, we go further with more hidden layers, fine-tuning parameters towards neural network model, and see how high accuracy can get.

Here (4.4) are learned features that resemble pen strokes. In other words, the model has learned to represent handwritten characters in terms of what pen strokes appear in an image.

The final model (4.5) we use has 2 hidden layers, and a final softmax classifier layer capable of classifying the MNIST digits as desired.

With complexity goes up, the perform of model depends highly on parameters. Without fine-tuning, the model gives an accuracy of 87.77%, but with fine-tuning, it gives an accuracy of 98.38%, which is similar to SVM with Gaussian kernel.

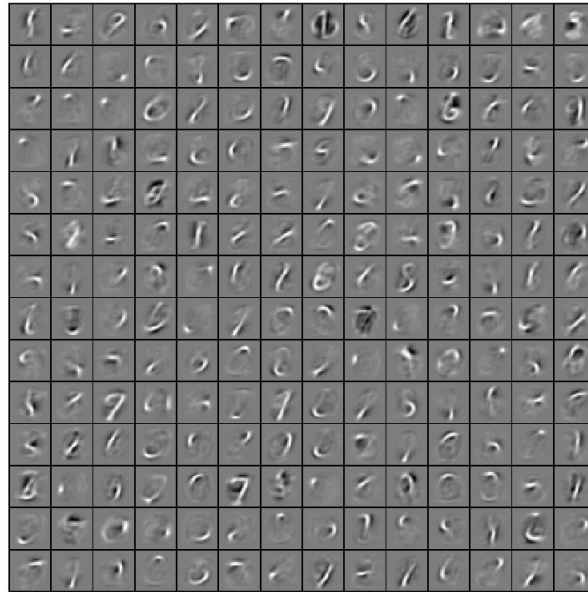


FIGURE 4.4: Learned Features of Handwritten Digits

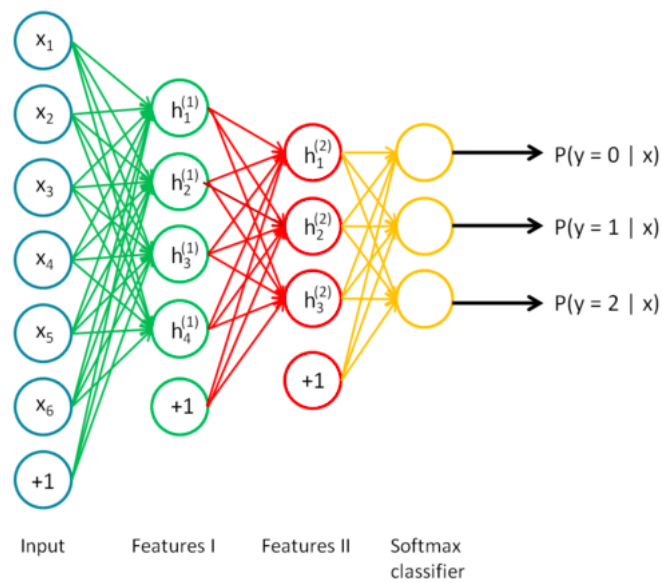


FIGURE 4.5: Deep Learning Model

Chapter 5

Conclusion

This paper focuses on spaces of functions generated by kernels, which are called reproducing kernel Hilbert spaces (RKHS). Upon this special space, we develop a kernel-based approach for classification in machine learning area. The remarkable feature of this method is while the criterion of classification problem is defined over an infinite-dimensional space, the solution can be finite-dimensional. With this notable feature, we can apply this method in large-scale nonlinear classification problems, which are tricky and time-consuming by ordinary classification algorithms.

The theoretical basis of RKHS lies on the research of Hilbert space in boundary value problems and reproducing property in the theory of integral equations. These spaces have wide applications, including complex analysis, harmonic analysis, and quantum mechanics. RKHS are particularly important in the field of statistical learning theory because of the celebrated Riesz representer theorem which states that every function in an RKHS can be written as a linear combination of the kernel function evaluated at the training points. Based on this property, we extend the theory into several widely-used statistical learning algorithms such as principle component analysis (PCA) and support vector machine (SVM), by the name of kernel methods, which owe their name to the use of kernel functions. Kernel methods use an approach called ‘kernel trick’ which enable them to operate in a high-dimensional, implicit feature space without ever computing the coordinates of the data in that space, but rather by simply computing the inner products between the images of all pairs of data in the feature space. This operation is often computationally cheaper than the explicit computation of the coordinates.

In our real-word example, we experiment on a large database *The MNIST Database of Handwritten Digits* (60,000 training, 10,000 testing observations), derived from standard NIST. These handwritten digits can never be standard with their sizes and most of them have some sort of rotation, all these special characteristics make feature selection difficult to apply. Our experiments have shown that the kernel methods perform extremely well with around 2 percent of errors, considering some handwritten digits are too hasty and careless to be labeled by humans, the result is really inspiring.

There are several problems which need to be addressed in further research on RKHS and its related algorithms. One of the challenges is the kernels used in algorithms must be positive-definite. This is reasonable since the theorem proves that given a positive definite kernel K , there is a unique associated RKHS with K as a reproducing kernel. Yet, positive-definite kernels are limited and some widely used kernels, such as sigmoid kernel, lack this property. If we can build up a method linking non-positive

definite kernels with RKHS (with some sort of limitation), we can apply this “kernel trick” into more algorithms and thus make it more powerful in real-world examples. Another difficulty lies in applying and evaluating algorithms. It’s difficult to make a distinction between several algorithms with almost the same testing accuracy shown at the end, thus make ourselves confusing to choose the most appropriate model, especially in unsupervised learning. More visualization and middle procedure should be displayed so that engineers can have more measurement towards these kernel-based methods.

Appendix A

Loss Function with Shrinkage

In this Appendix, we describe the formula for regularization problems.

A.1 Statistical Decision Theory

We seek a function $f(X)$ for predicting Y given values of the input X . This requires a loss function $L(Y, f(X))$ for penalizing errors in prediction, and by far the most common and convenient is squared error loss: $L(Y, f(X)) = (Y - f(X))^2$.

- Least squares assumes $f(X)$ is well approximated by a globally linear function,

$$f(X) = X^T \hat{\beta}$$

- k-nearest neighbors assumes $f(x)$ is well approximated by a locally constant function,

$$f(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$$

A.2 Shrinkage Methods

Regularization is designed to address the problem of overfitting.

High bias or underfitting is when the form of our hypothesis maps poorly to the trend of the data. It is usually caused by a function that is too simple or uses too few features.

At the other extreme, overfitting or high variance is caused by a hypothesis function that fits the available data but does not generalize well to predict new data. It is usually caused by a complicated function that creates a lot of unnecessary curves and angles unrelated to the data.

To address the issue of overfitting by regularization while still remaining all the features, we can keep all the features by reducing the parameters θ_j , this is exactly what shrinkage means. It works well when we have a lot of slightly useful features.

Using shrinkage methods, the cost function mentioned above can be shown as follows,

$$L(Y, f(X)) = \min_f \left\{ \sum_{i=1}^N (y_i - f(x_i))^2 + \lambda J(f) \right\}$$

where $\lambda \geq 0$ is a complexity parameter that controls the amount of shrinkage: the larger the value of λ , the greater the amount of shrinkage, and

thus cause higher bias and lower variance. The user-selected functional $J(f)$ will be large for functions f that vary too rapidly over small regions of input space.

A.3 Some Classical Regularization Problems

- Ridge Regression,

$$L(Y, f(X)) = \min_{\beta} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

- The Lasso,

$$L(Y, f(X)) = \min_{\beta} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

- Regularized Logistic Regression,

$$L(Y, f(X)) = \min_{\theta} \left\{ \left[\sum_{i=1}^N \log(h_{\theta}(x_i)) + (1 - y_i) \log(1 - h_{\theta}(x_i)) \right] + \lambda \sum_{j=1}^n \theta_j^2 \right\}$$

Appendix B

Gaussian Kernel

The Gaussian kernel $k(x, y) = \exp(-\frac{\|x-y\|^2}{2\sigma^2})$ is very popular and makes a good default kernel especially in absence of expert knowledge about data and domain. What makes the Gaussian kernel universal? How does the Gaussian kernel separate seemingly any sort of nonlinear data exceptionally well?

(B.1) is an example of feature space produced by Gaussian kernel. The feature space has infinite dimensions (maximum to the size of training sample), through adjusting the parameter σ , we can have a tradeoff between bias and variance. Small σ corresponds to low bias and high variance (overfitting), vise visa.

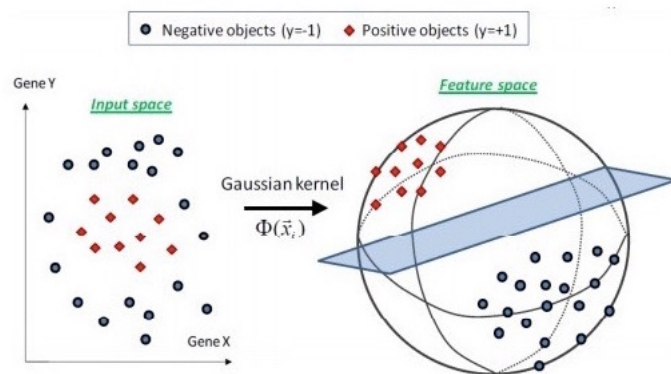


FIGURE B.1: Feature Space Produced by Gaussian Kernel

B.1 Locality

Given a bounded continuous function f , there exists a function $g \in \mathcal{H}$, such that f and g are close (in the sense of $\|\cdot\|_\infty$) up to arbitrary precision needed. Basically, this means Gaussian kernel gives functions which can approximate "nice" (bounded, continuous) functions arbitrarily well. Gaussian and Laplace kernels are universal. A polynomial kernel, for example, is not.

B.2 Separability

The Gaussian kernel for dimensions higher than one, say N , can be described as a regular product of N one-dimensional kernels. For example,

the Gaussian kernel be:

$$f(x; \sigma_1^2) = \frac{1}{\sqrt{2\pi}\sigma_1} e^{-x^2}$$

, then we have:

$$f(x, y; \sigma_1^2 + \sigma_2^2) = f(x; \sigma_1^2) f(y; \sigma_2^2).$$

B.3 Anisotropy

Anisotropy of a Gaussian kernel means that the scales, or standard deviations, are different for the different dimensions. When they are the same in all directions, the kernel is called isotropic, which means that the behavior of the function is in any direction the same. For 2D this means the Gaussian function is circular, for 3D it looks like a fuzzy sphere.

B.4 The Fourier transform of the Gaussian kernel

The basis functions of the Fourier transform \mathcal{F} are the sinusoidal functions $\exp(i\omega x)$. The definitions for the Fourier transform and its inverse are:

the Fourier transform:

$$F(\omega) = \mathcal{F}\{f(x)\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(x) e^{i\omega x} dx$$

the inverse Fourier transform:

$$\mathcal{F}^{-1}\{F(\omega)\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} F(\omega) e^{-i\omega x} d\omega$$

So the Fourier transform of the Gaussian function is again a Gaussian function, but now of the frequency ω . Also, the fact is smaller kernel in the spatial domain gives a wider kernel in the Fourier domain, and vice versa.

If we consider a Laplace kernel $k(x, y) = \exp(-\frac{\|x-y\|}{\sigma})$ which is also radial basis function kernel, its Fourier transform is a Cauchy distribution which drops much slower than the exponential function in the Fourier transform of a Gaussian kernel. This means a function f will have more high-frequency components. As a result, the function given by a Laplace kernel is "rougher" than that given by a Gaussian kernel.

Reference

- [1] Chih-Chung Chang and Chih-Jen Lin. “LIBSVM: A library for support vector machines”. In: *ACM Transactions on Intelligent Systems and Technology* 2 (3 2011). Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>, 27:1–27:27.
- [2] Kenji Fukumizu, Francis R. Bach, and Michael I. Jordan. “Dimensionality Reduction for Supervised Learning with Reproducing Kernel Hilbert Spaces”. In: *Journal of Machine Learning Research* 5 (2004), pp. 73–79.
- [3] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. 2nd. Springer-Verlag, 2009, pp. 167–173, 417–455, 547–550. URL: <http://statweb.stanford.edu/~tibs/ElemStatLearn/>.
- [4] Peter D. Lax. *Functional Analysis*. Wiley-Interscience, 2002, pp. 36–71.
- [5] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. *The MNIST Database of Handwritten Digits*. URL: <http://yann.lecun.com/exdb/mnist/>.
- [6] Andrew Ng. *Machine Learning*. URL: <http://openclassroom.stanford.edu/MainFolder/CoursePage.php?course=MachineLearning>.
- [7] Andrew Ng et al. *UFLDL Tutorial*. 2011. URL: http://deeplearning.stanford.edu/wiki/index.php/UFLDL_Tutorial.
- [8] Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. “Kernel Principal Component Analysis”. In: *Lecture Notes in Computer Science* 1327 (June 2006), pp. 583–588.
- [9] Grace Wahba. “Advances in Kernel Methods”. In: ed. by Bernhard Schölkopf, Christopher J. C. Burges, and Alexander J. Smola. MIT Press, 1999. Chap. Support Vector Machines, Reproducing Kernel Hilbert Spaces, and Randomized GACV, pp. 69–88.