

Spectral Analysis of Symmetric and Anti-Symmetric Pairwise Kernels

MARKUS VILJANEN

PRO GRADU-TUTKIELMA

TURUN YLIOPISTO

INFORMAATIOTEKNOLOGIAN LAITOS

TIETOJENKÄSITTELYTIEDE

MARRASKUU 2015

Contents

1	Prelude-----	1
1.1	Thesis topic-----	1
1.2	Thesis contribution -----	1
1.3	Research skills -----	2
2	Machine Learning-----	4
2.1	Conceptual framework-----	4
2.2	Approximation-----	5
2.3	Sampling -----	11
2.4	Learning-----	14
2.5	Hypotheses spaces -----	16
3	Reproducing Kernel Hilbert Space -----	25
3.1	Spaces associated to a kernel -----	25
3.2	Equivalence of spaces -----	27
3.3	Kernel and space modifications -----	31
3.4	Prior knowledge in Kernels -----	33
3.5	Examples of prior knowledge in Kernels-----	34
4	Theory of Linear Inverse Problems -----	37
4.1	Linear operators-----	37
4.2	Fundamental subspaces of a linear operator-----	38
4.3	Spectral Theorems -----	40
4.4	Singular Value Decompositions -----	41
4.5	Pseudoinverse of linear operators-----	43
4.6	Regularization of the pseudoinverse-----	45
5	Theory of Kernel Approximation -----	47
5.1	Random variable formulation -----	47

5.2	Abstract minimizers -----	49
5.3	Kernel minimizers -----	52
5.4	Spectral and singular value decompositions -----	54
5.5	Connection between \mathcal{H} and L^2 and Mercer's theorem-----	58
5.6	Convergence-----	61
6	Learning relations -----	68
6.1	Introduction to relations -----	68
6.2	Symmetric and anti-symmetric kernel -----	69
6.3	Examples of pairwise kernels -----	70
6.4	Approximation properties of symmetric and anti-symmetric kernels -----	72
6.5	Kernel matrices of symmetric and anti-symmetric kernels-----	74
6.6	Integral operators of symmetric and anti-symmetric kernels-----	78
6.7	Effective dimension of symmetric and anti-symmetric kernels-----	81
Appendix A: Effective dimension -----		1
A.1	Proof of $D\lambda SKS \leq D\lambda(K)$ -----	1
A.2	Effective dimension – a simple example-----	4
References -----		84

1 Prelude

1.1 Thesis topic

This thesis presents my small contribution to the topic of symmetric and anti-symmetric kernels in the context of a comprehensive theoretical framework developed by other researchers. The mathematical theory of kernels defining a function space saw its beginnings in the mid-20th century (Aronszajn, Theory of Reproducing Kernels, 1948), but the adaptation of kernels into machine learning gained significant traction only in the last two decades (Steinwart & Christmann, 2008). Mathematics formulating and investigating theoretical properties of this learning process is constantly evolving, and major learning theoretic results this thesis builds on are closer to only one decade old (Bauer; Pereverzev; & Rosasco, 2007).

The title of this thesis is inaccessible to a layman, but it can be described in a more general way as: "Suppose that we have a particular algorithm which is able to form a model of a quantifiable relationship after seeing only examples of it. If this algorithm is adapted to learn relations between things, and we know that those relations are symmetric or anti-symmetric, can we force it to learn only those relations and then guarantee that it learns better?" Given this question, a vague intuitive answer is a resounding yes. But mathematically formulating what this question exactly means, incorporating the solution into the algorithm and then answering the question is not trivial.

1.2 Thesis contribution

The research presented in the chapters up to the final one has been the work of other researchers, the aim was only to create a compendium of results required to understand the fundamental result of this thesis. This may have proven to be an ill-advised journey; the topic is very much more complicated than the author estimated. My small but central contribution in question is the result on the effective dimension of those kernels. After existence of this result, Prof. Pahikkala's productivity greatly outpaced mine to produce rest of the research draft (Pahikkala, Viljanen, Airola, & Waegeman, 2015).

Choosing this problem was my own decision from multiple options of Prof. Pahikkala's, and it was presented as it is stated. Is $D_\lambda(SKS) \leq D_\lambda(K)$? I did not even understand the question at first. Reviewing linear algebra and gaining significant intuition to the environment around the problem took approximately two months. Forming a big picture from Pahikkala's papers was a central part of this process. The actual solution used the intuition obtained to transform the problem into another formulation to which solution seemed geometrically obvious. A particular name existed for this formulation, principal submatrix, and a literature search arrived at Poincaré's result (Poincaré, 1890) from a different field which completed the final step. This 3-page original proof is contained in Appendix A.1.

An infinite-dimensional version of Poincaré's result was quite hard to track down, understanding the literature took almost a month and two weeks of following multiple references were required to arrive at a textbook which mentioned a similar result (Stenger & Alexander, 1972), but with formulation through complement names! Highly ironically the first legible statement was stated in (Aronszajn, Rayleigh-Ritz and A. Weinstein methods for approximation of eigenvalues: I. Operations in a Hilbert space, 1948), a paper by the originator of the monumental work on kernels (Aronszajn, Theory of Reproducing Kernels, 1948) published the very same year.

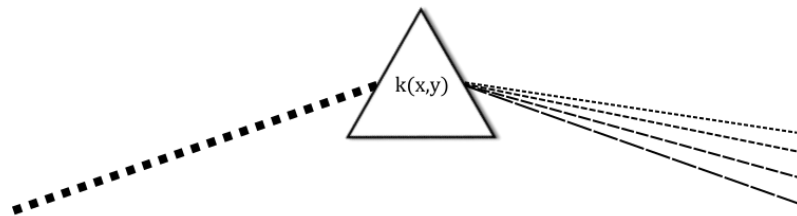
1.3 Research skills

From an academic perspective, the process of writing this thesis certainly cultivated research skills. The central piece of work was solving the problem, and this was the first time I was exposed to a difficult problem without any context. I did not know if the problem was solvable with the knowledge I had; if, when and how to acquire more knowledge. Not visible in this thesis are the attempts that did not go anywhere, or reviews of knowledge that were not fruitful for solving the problem. Those alone might total the length this thesis. Even though in hindsight much of that labour was wasted and naïve, developing these skills was a success because in the future I know how to be wiser about approaching similar problems.

Yet the proof of this result was a minor consumption of my time. Most of the time was spent orienting around the vast machine learning literature which was or only seemed to be related to the result. The task was not made easy by the fact that each

paper seemed to use a slightly different formulation and the required material was scattered among multiple papers. This is perhaps even more important research skill the process taught me; to plan and execute research despite vagueness and open-endedness, orient oneself around a vast body of proximally related literature and perform a sufficiently thorough yet efficient review.

There are many more research skills that I have in my stupor only recently realized. Much more optimal approach would have been to plan a smooth progression small goals instead of directly attempting sink-or-swim on a difficult problem. Collaboration was an essential part of research: wherever it seemed like there might be results about a problem I consulted the related literature. In the future I should also learn to better collaborate with other researchers, since orienting around the vastness of the research context may be possible only after a decade of experience.



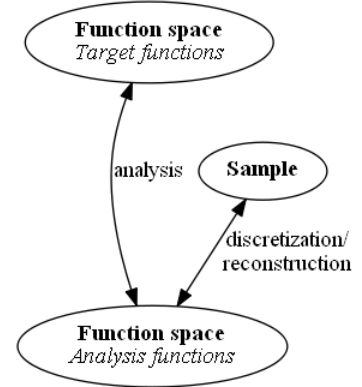
2 Machine Learning

2.1 Conceptual framework

The goal of supervised machine learning can be stated in a very simple way. We are provided data as $(\text{input}, \text{output})$ pairs:

INPUT	OUTPUT
(21.39, 15.52)	11.80
(18.17, 21.55)	0.59
(20.38, 16.29)	8.48
...	...

Mathematical framework

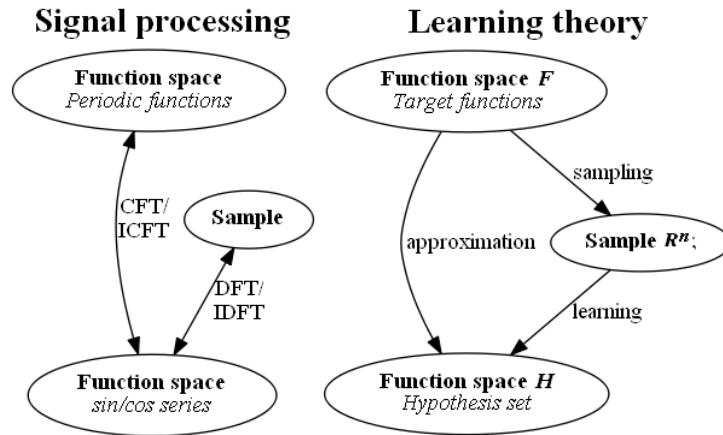


The goal is to form a model, such that for new `input`, the model gives an `output` which is close to the true unknown answer. Formally, we have:

- Input space X and output space Y
- Unknown target function $f: X \rightarrow Y$
- Data set $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$ where $y_i = f(x_i)$ and $x_i \in X, y_i \in Y$.
- A learning algorithm \mathcal{L} which returns a function $g = \mathcal{L}(D)$ such that $g \approx f$.
- The learning algorithm \mathcal{L} also often explicitly encompasses:
 - A hypothesis set H of functions $X \rightarrow Y$ from which the function is selected
 - An error measure $E(D, g)$ which defines the best approximation hypothesis.
 - An efficient algorithm to find $g \in H$ minimizing $E(D, g)$.

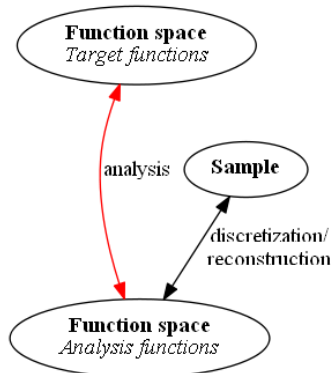
Machine learning is closely related to other fields, sharing mathematical foundations and methods in particular with learning theory, approximation theory, inverse problems, signal processing and statistics. The terminology is different however, and often the shared commonality is not obvious from the methods which make field-specific assumptions. For example, in signal processing, the Discrete Fourier transform (DFT) is used to transform a function into a discrete set of (point, value)-

pairs with a well defined inverse transform (IDFT). Given a sample satisfying certain requirements, this means we fully reconstruct an underlying periodic function from the sample. In machine learning, the goal is not to discretize and reconstruct known functions, and the lack of assumptions we can make about the target means full reconstruction is impossible. But both function-sample-transforms can be seen as instances of sampling.



This thesis aims to use the mathematical formalizations on the most general level possible, and then use their well-known special cases for practical applications. This allows a cross-disciplinary approach with increased concinnity at the cost of concreteness. The learning theory literature which encompasses RLS kernel methods and its learning bounds is often written in a highly general level as well, making varying special assumptions as needed.

2.2 Approximation



Linear combination of functions belonging to a certain function space can be used in the analysis of properties of a known target function or derive an expression for a function uniquely defined through some property. Function representation and approximation using a set of functions is a very important topic in mathematics, such as the use of power series in analysis, and has close connections to learning. Representation expresses the full function using the set of analysis functions, whereas approximation is an estimate with an error term of certain magnitude, possibly bounded locally. An accessible and application oriented resource is (Moon, 1999) whereas theoretical discussion is found for example in

(Kreyszig, 1989). One of the central mathematical tools used in approximation is the theory of abstract vector spaces and projections.

Projection Theorem (Kreyszig, 1989)

Let X be a Hilbert space and $X' \subseteq X$ a closed subspace of X . For any vector $\bar{x} \in X$, there exists a unique vector $\bar{x}_0 \in X'$ closest to \bar{x} : $\|\bar{x} - \bar{x}_0\| \leq \|\bar{x} - \bar{x}_1\|$ for all $\bar{x}_1 \in X'$. The point \bar{x}_0 is a minimizer of $\|\bar{x} - \bar{x}_0\|$ if and only if $\bar{x} - \bar{x}_0 \perp X'$.

2.2.1 Approximation with a finite set of vectors

Let X be a real Hilbert space and $X' = \text{span}\{p_1, \dots, p_m\}$ be a subspace spanned by linearly independent vectors p_1, \dots, p_m . Given a vector $\bar{x} \in X$, the approximation goal is to find vector $\hat{x} \in X'$ which approximates \bar{x} :

$$\bar{x} = \hat{x} + \bar{e} = \sum_{i=1}^m c_i \bar{p}_i + \bar{e}$$

The best approximation can be defined as the vector which minimizes the norm of the error:

$$\|\bar{x} - \hat{x}\| = \|\bar{e}\|$$

Using the projection theorem, the error is minimized when the error $\bar{e} = \bar{x} - \hat{x}$ is orthogonal to X' . This is the case if and only if the error is orthogonal to each vector that spans X' :

$$\left\langle \bar{x} - \sum_{i=1}^m c_i \bar{p}_i, \bar{p}_j \right\rangle = 0, j = 1, \dots, m$$

These conditions may be written as a matrix

$$\begin{pmatrix} \langle \bar{p}_1, \bar{p}_1 \rangle & \cdots & \langle \bar{p}_m, \bar{p}_1 \rangle \\ \vdots & \ddots & \vdots \\ \langle \bar{p}_1, \bar{p}_m \rangle & \cdots & \langle \bar{p}_m, \bar{p}_m \rangle \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{pmatrix} = \begin{pmatrix} \langle \bar{x}, \bar{p}_1 \rangle \\ \vdots \\ \langle \bar{x}, \bar{p}_m \rangle \end{pmatrix}$$

Define the Gramian matrix $G_{ij} = \langle \bar{p}_i, \bar{p}_j \rangle$ (Moon, 1999) and the expression can be written:

$$G\bar{c} = \bar{p}$$

The coefficients solved:

$$\bar{c} = G^{-1}\bar{p}$$

The Gramian is:

- Symmetric: follows from the symmetry of the inner product $\langle \bar{p}_i, \bar{p}_j \rangle = \langle \bar{p}_j, \bar{p}_i \rangle$.
- Positive semidefinite, furthermore positive definite iff p_1, \dots, p_m are linearly independent: $\bar{y}^* G \bar{y} = \sum_{i=1}^m \sum_{j=1}^m y_i y_j \langle \bar{p}_i, \bar{p}_j \rangle = \sum_{i=1}^m \sum_{j=1}^m \langle y_i \bar{p}_i, y_j \bar{p}_j \rangle = \langle \sum_{i=1}^m y_i \bar{p}_i, \sum_{j=1}^m y_j \bar{p}_j \rangle = \|\sum_{i=1}^m y_i \bar{p}_i\|^2 \geq 0$.

In the special case that the vectors p_1, \dots, p_m are orthogonal, we do not have to solve a linear system, because the coefficients are simply $c_j = \langle \bar{x}, \bar{p}_j \rangle / \langle \bar{p}_j, \bar{p}_j \rangle$. Furthermore, orthogonal vectors of unit length are called orthonormal and $c_i = \langle \bar{x}, \bar{p}_i \rangle$. Because of this simplification, in many approximation tasks the approximating set is defined so that the vectors are orthonormal.

2.2.2 Approximation with infinite set of basis vectors

Given a set of orthonormal vectors $\{u_1, \dots, u_n\}$, we may associate with any $x \in X$ the series

$$\hat{x} = \sum_{i=1}^n \langle \bar{x}, u_i \rangle u_i$$

Then from the considerations before, let $\bar{e} = \bar{x} - \hat{x}$ with the orthogonality $\bar{e} \perp \hat{x}$ and

$$\bar{x} = \hat{x} + \bar{e}$$

From the orthogonality property:

$$\|\bar{x}\|^2 = \|\hat{x}\|^2 + \|\bar{e}\|^2 \Rightarrow \|\bar{x}\| \geq \|\hat{x}\| \geq 0$$

For $n = 1, 2, \dots$, the sequence $\|\hat{x}\|^2 = \sum_{i=1}^n |\langle \bar{x}, u_i \rangle|^2$ is a monotonically increasing bounded sequence, which thus converges. This leads to the following.

Lemma: Bessel's inequality

Let (u_i) be a sequence of orthonormal vectors, then for any $x \in X$:

$$\sum_{i=1}^{\infty} |\langle x, u_i \rangle|^2 \leq \|x\|^2$$

Theorem: Converge of coefficients (Kreyszig, 1989)

The series $s_n = \sum_{i=1}^n \alpha_i u_i$ of partial sums of vectors converges in the norm of a Hilbert space X , i.e:

$$\|s_n - s\| \rightarrow 0 \text{ as } n \rightarrow \infty$$

if and only if the series $\alpha_n = \sum_{i=1}^n \alpha_i^2$ converges.

Proof

Because of orthonormality, $\|s_m - s_n\| = \|\alpha_m u_m + \dots + \alpha_{n+1} u_{n+1}\| = \alpha_m^2 + \dots + \alpha_{n+1}^2$, i.e. the series s_n is Cauchy iff α_n is Cauchy. Both X and \mathbb{R} are complete, in which series is Cauchy iff it converges.

In particular, this combined with Bessel's inequality implies that for any $\bar{x} \in X$, the series $\hat{x} = \sum_{i=1}^{\infty} \langle \bar{x}, u_i \rangle u_i$ converges to an element in X . The natural question is when $\hat{x} = \bar{x}$?

An essential condition is the totality of (u_i) , which can be shown to be equivalent to multiple conditions, of which Parseval's equality is common (Moon, 1999). The existence of countable total set (u_i) can be shown to be equivalent to separability of X (Kreyszig, 1989).

Definition

An orthonormal set $U = \{u_1, u_2 \dots\}$ in a Hilbert space X is total if its span is dense in X :

$$\overline{\text{span}(U)} = X$$

Lemma: Parseval's equality (Kreyszig, 1989).

An orthonormal set $U = \{u_1, u_2 \dots\}$ is total if and only if Bessel's inequality is an equality:

$$\forall \bar{x} \in X : \sum_{i=1}^{\infty} \langle \bar{x}, u_i \rangle^2 = \|\bar{x}\|^2$$

Furthermore, then $\bar{x} = \sum_{i=1}^{\infty} \langle \bar{x}, u_i \rangle u_i$.

2.2.3 Example: Linear regression as approximation in \mathbb{R}^n

Given a set of vectors $\{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n\}$ with $\bar{x}_i \in \mathbb{R}^d$ and outputs $\{y_1, \dots, y_n\}$, the goal of linear regression is to model an unknown linear relationship

$$x_{i1}w_1 + \dots + x_{id}w_d = y_i \text{ for } i = 1 \dots n$$

The best approximation is defined as the minimizer of the squared error $\sum_{i=1}^n (\sum_j x_{ij}w_j - y_i)^2$, which may be written in matrix form by stacking vectors as rows of $X = (\bar{x}_1, \dots, \bar{x}_n)^T$:

$$E(\bar{w}) = \|X\bar{w} - \bar{y}\|_2^2$$

But $\|\cdot\|_2^2$ is simply the Euclidean norm in the space \mathbb{R}^n and denoting the columns $X = (\bar{x}'_1, \dots, \bar{x}'_n)$ we have the approximator $X\bar{w} = \sum_{i=1}^n \bar{x}'_i w_i$. By the projection theorem the minimizer is the orthogonal projection in \mathbb{R}^n into the column space of X , therefore for the minimizer:

$$\begin{pmatrix} \bar{x}'_1 \bar{x}'_1 & \dots & \bar{x}'_n \bar{x}'_1 \\ \vdots & \ddots & \vdots \\ \bar{x}'_1 \bar{x}'_n & \dots & \bar{x}'_n \bar{x}'_n \end{pmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} = \begin{pmatrix} \bar{x}'_1 \bar{y} \\ \vdots \\ \bar{x}'_n \bar{y} \end{pmatrix}$$

Which may be written more succinctly in matrix form, since in this case $G = X^T X$:

$$X^T X \bar{w} = X^T \bar{y}$$

If the columns are linearly independent this can be solved directly.

2.2.4 Example: Approximation in $L^2[a,b]$

Let $L^2[0,1]$ be the Hilbert space obtained by taking the completion of the inner product space of continuous real valued functions on the interval $[0,1]$ with the inner product

$$\langle f, g \rangle = \int_0^1 f(t)g(t)dt$$

Suppose we want to find the best $m - 1$ degree polynomial approximation of a continuous function $f(t)$ in the sense that we minimize the norm induced by the inner product $\|f\| = \sqrt{\langle f, f \rangle}$:

$$\int_0^1 (f(t) - p(t))^2 dt$$

Take $1, t, \dots, t^{m-1}$ as the approximation vectors so that $p(t) = c_0 + c_1 t + c_2 t^2 + \dots + c_{m-1} t^{m-1}$. The minimum norm approximation problem becomes:

$$\begin{bmatrix} \langle 1, t \rangle & \cdots & \langle 1, t^{m-1} \rangle \\ \vdots & \ddots & \vdots \\ \langle t^{m-1}, 1 \rangle & \cdots & \langle t^{m-1}, t^{m-1} \rangle \end{bmatrix} \begin{bmatrix} c_0 \\ \vdots \\ c_m \end{bmatrix} = \begin{bmatrix} \langle f, 1 \rangle \\ \vdots \\ \langle f, t^{m-1} \rangle \end{bmatrix}$$

Where:

$$\langle t^i, t^j \rangle = \int_0^1 t^{i+j} dt = \frac{1}{i+j+1}$$

The Gramian is the well known ill-conditioned Hilbert matrix:

$$H = \begin{bmatrix} 1 & \cdots & \frac{1}{m} \\ \vdots & \ddots & \vdots \\ \frac{1}{m} & \cdots & \frac{1}{2m} \end{bmatrix}$$

Suppose we solved the coefficients (c_i) of the above problem $H\bar{c} = \bar{f}$ through taking the inverse $\bar{c} = H^{-1}\bar{f}$. The solution then is:

$$p(t) = c_0 + c_1 t + c_2 t^2 + \cdots + c_{m-1} t^{m-1}$$

In practice we might use the Legendre polynomials, which are obtained by applying the Gram-Schmidt algorithm to $\{1, t, t^2 \dots\}$ and also form a total orthonormal set on $L^2[0,1]$ (Kreyszig, 1989).

2.2.5 Example: Fourier series

Any (complex) function periodic on $[0, 2\pi)$ may be represented using the series (Moon, 1999):

$$f(t) = \sum_{n=-\infty}^{\infty} \langle x, p_n \rangle p_n(t)$$

Where $p_n(t)$ are the following orthonormal Fourier basis functions:

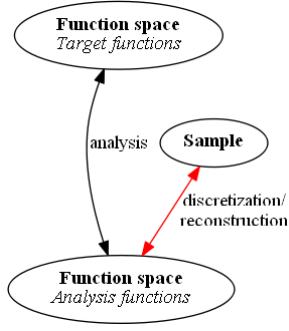
$$p_n(t) = \frac{1}{\sqrt{2\pi}} e^{nti}$$

The coefficients can be found simply by:

$$\langle x, p_n \rangle = \frac{1}{\sqrt{2\pi}} \int_0^{2\pi} f(t) e^{-nti} dt$$

The representation is the Fourier transform of $f(t)$.

2.3 Sampling



Sampling a function is an important mathematical concept, but the context is very different in machine learning compared to those often occurring in theoretical discussion or applications such as signal processing. This setting has been recently used in different forms in the machine learning literature, see for example the discussion in (De Vito; Rosasco; Caponetto; De Giovannini; & Odone, 2005).

Connection to learning theory and clear formalization were recently given in (Smale & Zhou, Shannon Sampling and Function reconstruction, 2004) and (Smale & Zhou, Shannon sampling II: Connections to learning theory, 2005). Below we present their formalization and relate it to our goal.

2.3.1 Discretization of a function

Let \mathcal{H} be a space of functions of the form $f: X \rightarrow \mathbb{R}$. In discretization, the goal is to define $(x_1, \dots, x_n) \in X^n$ such that the mapping $\mathcal{F}_{\bar{x}}: f \mapsto (f(x_1), \dots, f(x_n))$ is a **bijection** $\mathcal{F}_{\bar{x}}: \mathcal{H} \rightarrow \mathbb{R}^n$. This is unlike sampling in machine learning, where the sample is not purposefully chosen and it often is impossible to fully reconstruct the function from the sample, in other words the inverse $\mathcal{F}_{\bar{x}}^{-1}$ does not exist. If the inverse does exist, the map $\mathcal{F}_{\bar{x}}$ is said to discretize the function space. The two following examples of discretization are widely used, see for example (Smale & Zhou, Shannon sampling II: Connections to learning theory, 2005).

Example 1

Consider a space of polynomials \mathcal{H}_d up to degree d : $p(x) = a_0 + a_1x + \dots + a_dx^d$. Given any sample of $d + 1$ points $\bar{x} = (x_1, \dots, x_{d+1})$ its value may be computed:

$$\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_{d+1}) \end{pmatrix} = \begin{pmatrix} 1 & \dots & x_1^d \\ \vdots & \ddots & \vdots \\ 1 & \dots & x_{d+1}^d \end{pmatrix} \begin{pmatrix} a_0 \\ \vdots \\ a_d \end{pmatrix}$$

More succinctly, we have a linear transformation $L_{\bar{x}}: \mathbb{R}^{d+1} \rightarrow \mathbb{R}^{d+1}$ known as the Vandermonde matrix:

$$f(\bar{x}) = L_{\bar{x}} \bar{a}$$

When the x_i are distinct $L_{\bar{x}}: \bar{a} \mapsto (f(x_i))_{i=1}^{d+1}$ is a bijection, since the determinant is nonzero. Similarly the map $\sum_{i=0}^d a_i x^i \rightarrow \bar{a}$ from coefficients to the polynomial is a bijection, so their composition is a discretizing bijection $\mathcal{F}_{\bar{x}}: \mathcal{H}_d \rightarrow \mathbb{R}^{d+1}$.

Example 2 (Whittaker-Shannon-Nyquist Sampling Theorem)

Let $\phi(x) = \frac{\sin(\pi x)}{\pi x}$ and $\phi_t(x) = \phi(x - t)$. If a function $f \in L^2(\mathbb{R})$ has its Fourier transform supported in $[-\pi, \pi]$, then

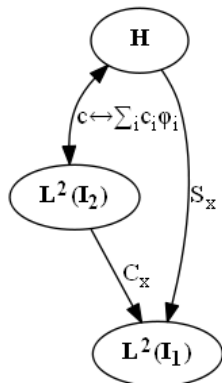
$$f = \sum_{t \in \mathbb{Z}} f(t) \phi_t$$

Functions $\{\sum_{t \in \mathbb{Z}} a_t \phi_t : \bar{a} \in \mathbb{R}^{\mathbb{Z}}\}$ are the representation space, and as a consequence of the theorem the bijection $\mathcal{F}_{\bar{x}}: f \mapsto (f(t))_{t \in \mathbb{Z}}$ exists.

2.3.2 Reconstruction of a function

Following the formulation of (Smale & Zhou, Shannon Sampling and Function reconstruction, 2004), we present the sampling framework for machine learning as follows. Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be function and for $x_i \in \mathbb{R}^d$ denote the finite or infinite set where the function has been sampled $\bar{x} = \{x_1, x_2, \dots\} = \{x_i\}_{i \in I_1}$. Assuming countability the index set I_1 is $\{1, \dots, n\}$ or \mathbb{N} . Define the sample space $\ell^2(I_1)$ ($= \mathbb{R}^n$ or ℓ^2), and denote the values of f at \bar{x} as $\bar{y} = (y_i)_{i \in I_1} \in \ell^2(I_1)$.

Likewise, let \mathcal{H} be a Hilbert space of continuous functions and let $\{\phi_i\}_{i \in I_2}$ be a finite or countably infinite dimensional orthonormal basis of \mathcal{H} . Define the coefficient space $\ell^2(I_2)$ as the set of square summable sequences $(a_i)_{i \in I_2}$ with the inner product $\langle a, b \rangle = \sum_{i \in I_2} a_i b_i$. Then the mapping $\sum_{i \in I_2} c_i \phi_i \mapsto \bar{c}$ between $\mathcal{H} \rightarrow \ell^2(I_2)$ is an isomorphism. The space of sequences is defined so that instead of $f \in \mathcal{H}$ we may equivalently consider the coefficient sequence $(a_i)_{i \in I_2}$ of the series $f = \sum_{i \in I_2} a_i \phi_i$:



The sampling operator $S_{\bar{x}}$ is a map from the function space $\mathcal{H} \rightarrow \ell^2(I_1)$:

$$S_{\bar{x}}(f) = (f(x_i))_{i \in I_1}$$

The sampling operator from the sequence space $\mathcal{C}_{\bar{x}}: \ell^2(I_2) \rightarrow \ell^2(I_1)$ is an infinite matrix:

$$\mathcal{C}_{\bar{x}} = (\phi_j(x_i))_{i \in I_1, j \in I_2}$$

This is well defined because it denotes the linear operator

$$(\mathcal{C}_{\bar{x}}\bar{c})_i = \sum_{j \in I_2} c_j \phi_j(x_i)$$

Considering either $\mathcal{C}_{\bar{x}}$ or $S_{\bar{x}}$ as the sampling operator, there are two fundamental problems with its inverse. Since the sample was provided and not purposefully chosen, the operator may be neither injective nor surjective. If the inverse is not a bijection, we need to clarify how we define function reconstruction.

- a) Surjectivity: If the function space \mathcal{H} is not expressive enough on a large enough sample \bar{x} , for values $(y_i)_{i \in I_1}$ there may not be a function f such that $S_{\bar{x}}f = \bar{y}$.
- b) Injectivity: If the function space \mathcal{H} is expressive enough or the sample \bar{x} too small, there may very well be two functions f_1, f_2 such that $S_{\bar{x}}f_1 = S_{\bar{x}}f_2$.

Because it may be impossible to use an inverse operator, instead of full function reconstruction one needs to formulate a well-defined minimization problem. Sometimes injectivity or surjectivity may be assumed, or the existence problem of the inverse solved by:

- a) Defining an approximation task, such as minimizing an error $\mathcal{E}_{\bar{x}, \bar{y}}(f)$.
- b) Defining a property which makes the solution $f \in \mathcal{H}$ unique, such as penalized or minimum norm.

Aside from injectivity and surjectivity, there may be an additional problem with the stability of the solution. It will be shown that in the minimum norm or unsubstantially penalized solution reconstructed from a sample of an unknown function with little noise added produces a very different solution. In learning this expressed as generalization ability, which this solution does not have: The learned function \hat{f} approximates the target f in the known sample \bar{x} very closely: $S_{\bar{x}}\hat{f} \cong S_{\bar{x}}f$. But given a new sample \bar{x}' , the values $S_{\bar{x}'}\hat{f}$ it produces can be wildly different from $S_{\bar{x}'}f$.

Clearly new tools are needed to determine a suitable hypothesis set \mathcal{H} and the rule which picks an approximating function $\hat{f} \in \mathcal{H} \approx f$ based on the sample $S_{\bar{x}}f$. Aside from injectivity, surjectivity and generalization, these tools need to incorporate considerations of the computational efficiency as well. Because there was nothing sampling operator specific in the formulation, for full generality it may be studied as an abstract linear inverse problem.

2.4 Learning

2.4.1 Learning as an inverse problem

Definition: Hadamard's criteria

Consider an abstract linear operator $T: \mathcal{X} \rightarrow \mathcal{Y}$ equation of the form

$$Tx = y$$

Hadamard's definition of a well-posed inverse problem, with equivalent conditions for the special case of a linear operator in parenthesis, is (Engl; Hanke; & Neubauer, 1996):

- a) For all admissible data, a solution exists ($\mathcal{R}(T) = \mathcal{Y}$)
- b) For all admissible data, a solution is unique ($\mathcal{N}(T) = \{0\}$)
- c) The solution depends continuously on the data (T^{-1} is continuous/bounded)

Supervised learning is an algorithm which solves a well-posed inverse problem. Since $S_{\bar{x}}: \mathcal{H} \rightarrow \ell^2(I_1)$ does not necessarily satisfy these conditions, a new problem formulation is needed.

2.4.2 Least squares and regularized least squares

In this thesis we focus on one particular learning algorithm, called **regularized least squares (RLS)**, which is widely used, effective and furthermore has very useful analytical and statistical properties (Cucker & Smale, 2001). It solves this problem as follows:

- a) The goal is to minimize the least squares error,:

$$\mathcal{E}_{\bar{x}, \bar{y}}(f) = \sum_{i \in I_1} (f(x_i) - y_i)^2$$

- b) Define the optimal solution as the solution with minimum norm $|||_{\mathcal{H}}$.
- c) Possibly regularize (penalize) the norm to achieve desirable stability properties.

For example in (Smale & Zhou, Shannon Sampling and Function reconstruction, 2004), it is assumed that \bar{x} provides rich data with respect to \mathcal{H} in the sense that

$\inf_{v \in \ell^2(I_2)} \frac{\|\mathcal{C}_{\bar{x}} v\|_{\ell^2(I_1)}}{\|v\|_{\ell^2(I_2)}} > 0$, which is equivalent to $\mathcal{C}_{\bar{x}}^T \mathcal{C}_{\bar{x}}$ having bounded inverse which means that the map is injective. The solution is then easily obtained:

Theorem: Least Squares solution (Smale & Zhou, Shannon Sampling and Function reconstruction, 2004)

The problem is to reconstruct f^* from $\bar{y} \in \ell^2(I_1)$. Consider the minimization problem

$$\operatorname{argmin}_{f \in \mathcal{H}} \sum_{i \in I_1} (f(x_i) - y_i)^2$$

If \bar{x} provides rich data, the solution is

$$f_{\bar{x}} = \sum_{i \in I_2} c_i \phi_i$$

Where $\bar{c} = L\bar{y}$ with $L = (\mathcal{C}_{\bar{x}}^T \mathcal{C}_{\bar{x}})^{-1} \mathcal{C}_{\bar{x}}^T$, which is the Moore-Penrose inverse of $\mathcal{C}_{\bar{x}}$.

Proof

This is a quadratic form that may be expressed:

$$Q(\bar{c}) := \|\mathcal{C}_{\bar{x}, \bar{t}} \bar{c} - \bar{y}\|_{\ell^2(I_2)}^2 = \langle \mathcal{C}_{\bar{x}}^T \mathcal{C}_{\bar{x}} \bar{c}, \bar{c} \rangle_{\ell^2(I_2)} - 2 \langle \mathcal{C}_{\bar{x}} \bar{c}, \bar{y} \rangle_{\ell^2(I_1)} + \langle \bar{y}, \bar{y} \rangle_{\ell^2(I_1)}$$

Take the functional derivative and equate $\frac{\delta}{\delta \bar{c}} Q(\bar{c}) = 0$ to find the minimizer:

$$\begin{aligned} \mathcal{C}_{\bar{x}}^T \mathcal{C}_{\bar{x}} \bar{c} &= \mathcal{C}_{\bar{x}}^T \bar{y} \\ \bar{c} &= (\mathcal{C}_{\bar{x}}^T \mathcal{C}_{\bar{x}})^{-1} \mathcal{C}_{\bar{x}}^T \bar{y} \end{aligned}$$

If the function space is a Hilbert space, one possible choice for the property which removes the requirement for the injectivity assumption and makes the solution unique is minimum norm or penalized norm solution. To use a penalized norm, define a new error, for $\lambda > 0$:

$$\mathcal{E}_{\bar{x}, \bar{y}}(f) = \operatorname{argmin}_{f \in \mathcal{H}} \sum_{x \in I_1} (f(x_i) - y_i)^2 + \lambda \|f\|_H$$

Theorem: Regularized Least Squares solution (Smale & Zhou, Shannon Sampling and Function reconstruction, 2004)

The problem is to reconstruct f^* from $\bar{y} \in \ell^2(I_1)$. Consider the minimization problem

$$\operatorname{argmin}_{f \in \mathcal{H}} \sum_{x \in I_1} (f(x_i) - y_i)^2 + \lambda \|f\|_H$$

For an orthonormal basis $\{\phi_i\}_{i \in I_2}$ of \mathcal{H} , the solution is

$$f_{\bar{x}} = \sum_{i \in I_2} c_i \phi_i$$

Where $\bar{c} = L\bar{y}$ with $L = (\mathcal{C}_{\bar{x}}^T \mathcal{C}_{\bar{x}} + \lambda I)^{-1} \mathcal{C}_{\bar{x}}^T$, which is the Moore-Penrose inverse of $\mathcal{C}_{\bar{x}}$.

Proof

This is a quadratic form that may be expressed:

$$\begin{aligned} Q(\bar{c}) &:= \|\mathcal{C}_{\bar{x}} \bar{c} - \bar{y}\|_{\ell^2(I_1)}^2 + \lambda \|\bar{c}\|_{\ell^2(I_2)}^2 \\ &= \langle \mathcal{C}_{\bar{x}}^T \mathcal{C}_{\bar{x}} \bar{c}, \bar{c} \rangle_{\ell^2(I_2)} - 2 \langle \mathcal{C}_{\bar{x}} \bar{c}, \bar{y} \rangle_{\ell^2(I_1)} + \langle \bar{y}, \bar{y} \rangle_{\ell^2(I_1)} + \lambda \langle \bar{c}, \bar{c} \rangle_{\ell^2(I_2)} \end{aligned}$$

Take the functional derivative and equate $\frac{\delta}{\delta \bar{c}} Q(\bar{c}) = 0$ to find the minimizer:

$$\begin{aligned} (\mathcal{C}_{\bar{x}}^T \mathcal{C}_{\bar{x}} + \lambda I) \bar{c} &= \mathcal{C}_{\bar{x}}^T \bar{y} \\ \bar{c} &= (\mathcal{C}_{\bar{x}}^T \mathcal{C}_{\bar{x}} + \lambda I)^{-1} \mathcal{C}_{\bar{x}}^T \bar{y} \end{aligned}$$

2.5 Hypotheses spaces

In this section we consider different hypothesis spaces \mathcal{H} from which the sampling operator maps to the sample space $S_{\bar{x}}: \mathcal{H} \rightarrow \ell^2(I_1)$. As before, denote \mathcal{H} a Hilbert space of continuous functions and $\{\phi_i\}_{i \in I_2}$ an orthonormal basis of \mathcal{H} , so that $\overline{\operatorname{span}(\{\phi_i\}_{i \in I_2})}$ forms an analysis space. Because of the isomorphism, we may then equivalently operate on the coefficient sequence $(a_i)_{i \in I_2} \in \ell^2(I_2)$ of the series $f = \sum_{i \in I_2} a_i \phi_i$. Given a sample $\bar{y} \in \ell^2(I_1)$, to solve the problem one forms the sampling operator $\mathcal{C}_{\bar{x}}: \ell^2(I_2) \rightarrow \ell^2(I_1)$:

$$\mathcal{C}_{\bar{x}} = \left(\phi_j(x_i) \right)_{x \in I_1, j \in I_2}$$

The minimum norm closest $\|\cdot\|_2$ approximation solution is given by the pseudoinverse (see chapter 4):

$$\bar{c}^\dagger = \mathcal{C}_{\bar{x}}^\dagger \bar{y}$$

Regularized solution is essential if \mathcal{H} is expressive or sample noisy, Tikhonov regularization gives:

$$\bar{c}_\lambda^\dagger = (\mathcal{C}_\bar{x}^T \mathcal{C}_\bar{x} + \lambda I)^{-1} \mathcal{C}_\bar{x}^T \bar{y}$$

Solving for the coefficients \bar{c} , the solution is then given by the isomorphism between $\ell^2(I_2)$ and \mathcal{H} :

$$\bar{c} \mapsto \sum_{i \in I_2} c_i \phi_i$$

2.5.1 Hypothesis spaces

Linear hypotheses

Previous examples used functions $f: \mathbb{R} \rightarrow \mathbb{R}$ with 1-dimensional domain. In general, the target function may be of the form $f: \mathbb{R}^d \rightarrow \mathbb{R}$. Consider now the linear hypothesis space \mathcal{H}_d of functions $\mathbb{R}^d \rightarrow \mathbb{R}$:

$$\mathcal{H}_d = \{f(\bar{x}) = \bar{w}^T \bar{x} \mid \bar{w} \in \mathbb{R}^d\}$$

Since by definition for any function $f(\bar{x}) \in \mathcal{H}_d$:

$$f(\bar{x}) = w_1 x_1 + \dots + w_d x_d$$

A basis for this function space is the set of functions $\{\phi_i\}_{i=1}^d$ where $\phi_i(\bar{x}) = \bar{x}_i$, i.e the set $\{x_1, \dots, x_d\}$, with the coefficient space $\bar{w} \in \mathbb{R}^d$. For samples $\{\bar{x}_1, \dots, \bar{x}_n\}$, the sampling matrix is:

$$\mathcal{C}_\bar{x} = \begin{pmatrix} (\bar{x}_1)_1 & \dots & (\bar{x}_1)_d \\ \vdots & \ddots & \vdots \\ (\bar{x}_n)_1 & \dots & (\bar{x}_n)_d \end{pmatrix}$$

Given a set of samples $\{\bar{x}_1, \dots, \bar{x}_n\}$ and values $\{y_1, \dots, y_n\}$, the solution is obtained with previous tools.

Polynomial hypotheses

Linear hypotheses are very simple. The polynomial hypothesis up to degree m may be generalized from $\mathbb{R} \rightarrow \mathbb{R}$ to $\mathbb{R}^d \rightarrow \mathbb{R}$. Denote the polynomial hypothesis space \mathcal{H}_d^m of m -degree polynomials $\mathbb{R}^d \rightarrow \mathbb{R}$. Then we may define the following set of integer vectors whose entries sum to less than or equal to m :

$$\Omega = \{\bar{n} \in \mathbb{N}^d \mid \sum n_i \leq m\}$$

The space may then be expressed concisely:

$$\mathcal{H}_d = \left\{ \sum_{\sigma \in \Omega} c_{\sigma} x_1^{\sigma_1} \dots x_d^{\sigma_d} \mid \bar{c} \in \mathbb{R}^{|\Omega|} \right\}$$

In other words, every $f \in \mathcal{H}_d$ may be written:

$$f(\bar{x}) = c_1 x_1^m + c_2 x_1^{m-1} + \dots + c_{m+1} x_1^{m-1} x_2 + \dots + c_{|\Omega|} x_d^m$$

Every function is a linear function of the monomials $\phi_{\sigma}(\bar{x}) = x_1^{\sigma_1} \dots x_d^{\sigma_d}$ where the sum of the exponents is equal to or less than m . Therefore, the monomials $\{x_1^{\sigma_1} \dots x_d^{\sigma_d}\}_{\sigma \in \Omega}$ form a basis of \mathcal{H}_d with the coefficient space $\bar{c} \in \mathbb{R}^{|\Omega|}$. For samples $\{\bar{x}_1, \dots, \bar{x}_n\}$, the sampling matrix is:

$$\mathcal{C}_{\bar{x}} = \begin{pmatrix} \phi_{\sigma_1}(\bar{x}_1) & \dots & \phi_{\sigma_{|\Omega|}}(\bar{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_{\sigma_1}(\bar{x}_n) & \dots & \phi_{\sigma_{|\Omega|}}(\bar{x}_n) \end{pmatrix}$$

Given a set of samples $\{\bar{x}_1, \dots, \bar{x}_n\}$ and values $\{y_1, \dots, y_n\}$, the solution is obtain with previous tools. However, the number of these monomials grows rapidly as a function of the degree m . The $n \times |\Omega|$ matrix $\mathcal{C}_{\bar{x}}$ almost certainly has multiple solutions, a high degree of instability and the inverse lends itself to growing computational complexity.

Abstract feature space hypotheses

A tool often used in machine learning is the feature transform, let $\phi(\bar{x}) = \bar{u}$ be an arbitrary mapping $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^z$ to a high dimensional space. Define the hypothesis space \mathcal{H}_d^z as the set of linear hypothesis $\mathbb{R}^z \rightarrow \mathbb{R}$, which may correspond to a highly complicated nonlinear hypothesis $\mathbb{R}^d \rightarrow \mathbb{R}$:

$$\mathcal{H}_d^z = \{\bar{w}^T \phi(\bar{x}) \mid \bar{w} \in \mathbb{R}^z\}$$

An obvious basis for this function space is the set of functions $\{\phi(\bar{x})_i\}_{i=1}^z$, consisting of each entry of the feature vector $\phi(\bar{x})_i$ with the coefficient space $(w_i)_{i=1}^z \in \mathbb{R}^z$. For samples $\{\bar{x}_1, \dots, \bar{x}_n\}$, the sampling matrix can be written as the feature vectors as rows:

$$\mathcal{C}_{\bar{x}} = \begin{pmatrix} \phi(\bar{x}_1)^T \\ \vdots \\ \phi(\bar{x}_n)^T \end{pmatrix}$$

Depending on the dimensionality of the feature transform $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^z$, the $n \times z$ matrix $\mathcal{C}_{\bar{x}}$ may be very large which means the problem possibly has multiple solutions, instability and growing computational complexity.

2.5.2 Kernel from a feature space

The problem with the feature transform is that the dimension of \mathbb{R}^z may be large, thus explicitly performing this transform for each point and solving the linear system could make the approach computationally prohibitive. Deeper investigation into this setting reveals a trick one can utilize to overcome computational problems.

Let $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^z$ be an arbitrary feature transform. Define \mathcal{H}_d^z as the set of linear hypothesis in $\mathbb{R}^z \rightarrow \mathbb{R}$:

$$\mathcal{H}_d^z = \{\langle \bar{w}, \phi(\bar{x}) \rangle \mid \bar{w} \in \mathbb{R}^z\}$$

Where we have replaced the Euclidean inner product $\bar{w}^T \phi(\bar{x})$ by the general form. Suppose we have a set of basis functions for $f \in \mathcal{H}_d^z$:

$$f = c_1 \langle \phi(\bar{x}), \bar{w}_1 \rangle + \dots + c_z \langle \phi(\bar{x}), \bar{w}_z \rangle$$

The approximation problem is to solve, denoting the sampling matrix by Φ^T :

$$\begin{pmatrix} \hat{f}(x_1) \\ \vdots \\ \hat{f}(x_n) \end{pmatrix} = \begin{pmatrix} \langle \phi(\bar{x}_1), \bar{w}_1 \rangle & \dots & \langle \phi(\bar{x}_1), \bar{w}_z \rangle \\ \vdots & \ddots & \vdots \\ \langle \phi(\bar{x}_n), \bar{w}_1 \rangle & \dots & \langle \phi(\bar{x}_n), \bar{w}_z \rangle \end{pmatrix} \begin{pmatrix} c_1 \\ \vdots \\ c_z \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_z \end{pmatrix}$$

In the previous setting we used the obvious choice $\bar{w}_i = \bar{e}_i$, but then Φ^T is an $n \times z$ matrix and the problem is to solve $\Phi^T \bar{c} = \bar{y}$. Instead, we make the claim that we can select $\bar{w}_i = \phi(\bar{x}_i)$. Then we have an $n \times n$ matrix G_ϕ and the problem becomes:

$$G_\phi \bar{c} = \bar{y}$$

Theorem (adapted from (Moon, 1999)):

Let $\bar{\phi}_1, \dots, \bar{\phi}_n \in \mathbb{R}^z$ be linearly independent and let $V = \text{span}\{\bar{\phi}_1, \dots, \bar{\phi}_n\}$. Then $\|\bar{w}\| \in \mathbb{R}^z$ with minimum norm satisfying:

$$\begin{aligned} \langle \bar{w}, \bar{\phi}_1 \rangle &= y_1 \\ &\vdots \\ \langle \bar{w}, \bar{\phi}_n \rangle &= y_n \end{aligned}$$

belongs to V : $\bar{w} = \sum_{i=1}^n w_i \bar{\phi}_i$ with the coefficients \bar{w} satisfying:

$$\begin{pmatrix} \langle \bar{\phi}_1, \bar{\phi}_1 \rangle & \cdots & \langle \bar{\phi}_n, \bar{\phi}_1 \rangle \\ \vdots & \ddots & \vdots \\ \langle \bar{\phi}_1, \bar{\phi}_n \rangle & \cdots & \langle \bar{\phi}_n, \bar{\phi}_n \rangle \end{pmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

Proof:

Suppose \bar{w}_0 is a solution. Then so is any $\bar{w}_0 + \bar{n}$, where $\bar{n} \perp \bar{\phi}_1, \dots, \bar{\phi}_n$, i.e. $\bar{n} \in V^\perp$. In fact any solution \bar{w} can be written $\bar{w} = \bar{w}_0 + \bar{n}$ with $\bar{n} \in V^\perp$ (\bar{w}_1 also a solution $\Rightarrow \Phi^T(\bar{w}_1 - \bar{w}_0) = 0 \Rightarrow \bar{w}_1 - \bar{w}_0 \in V^\perp$)).

The problem is to find $\operatorname{argmin}_{\bar{n} \in V^\perp} \|\bar{w}_0 + \bar{n}\|$. This is equivalent to $\operatorname{argmin}_{\bar{n} \in V^\perp} \|\bar{n} - (-\bar{w}_0)\|$,

finding the minimum norm approximation of $-\bar{w}_0$ in the closed subspace V^\perp . The projection theorem states that the error of optimal solution \bar{n}^* is orthogonal to V^\perp : $\bar{n}^* - (-\bar{w}_0) \in (V^\perp)^\perp = V$. This proves that the optimal solution $\bar{w}^* = \bar{w}_0 + \bar{n}^* \in V$. Denote $\bar{w} = \sum_{i=1}^n w_i \bar{\phi}_i$ and take the inner products with $\bar{\phi}_1, \dots, \bar{\phi}_n$ to obtain the efficient formulation.

One implication of this is that instead of the potentially massive $n \times z$ system, we can solve $n \times n$ system G_ϕ . Furthermore, since the solution only depends on the dot products between feature transforms, even the feature transformation does not have to be performed if we have a **kernel** function k such that

$$k(\bar{x}_i, \bar{x}_j) = \langle \phi(\bar{x}_i), \phi(\bar{x}_j) \rangle$$

Utilizing $k(\bar{x}_i, \bar{x}_j)$ as a black box means that we have an algorithm which may have computationally better qualities but is equivalent to one with the feature transform. The underlying feature transform is implicit. The problem then is expressed through the **kernel matrix**:

$$\begin{pmatrix} k(\bar{x}_1, \bar{x}_1) & \cdots & k(\bar{x}_n, \bar{x}_1) \\ \vdots & \ddots & \vdots \\ k(\bar{x}_1, \bar{x}_n) & \cdots & k(\bar{x}_n, \bar{x}_n) \end{pmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

With the solution

$$f(\bar{x}) = w_1 k(\bar{x}, \bar{x}_1) + \cdots + w_n k(\bar{x}, \bar{x}_n)$$

2.5.3 Kernel from a function space

Suppose we have a Hilbert space \mathcal{H} of functions $f: X \rightarrow \mathbb{R}$ with a total orthonormal set $\{\phi_i\}_{i \in I_2}$ forming the analysis space, with the evaluation functional $\delta_x: f \mapsto f(x)$ continuous. Then for every $f \in \mathcal{H}$ there exists a representation $f = \sum_{i \in I_2} a_i \phi_i$ and a corresponding sequence $(a_i)_{i \in I_2} \in \ell^2(I_2)$.

Denoting the space of functionals \mathcal{H}' , by the Riesz representer theorem (Kreyszig, 1989, s. 188), for every bounded functional $\mathcal{g} \in \mathcal{H}'$ there exists a representer $g \in \mathcal{H}$ such that: $\mathcal{g}(f) = \langle f, g \rangle_{\mathcal{H}} \forall f \in \mathcal{H}$. Denote this map $\mathcal{R}: \mathcal{H}' \rightarrow \mathcal{H}$. Since evaluation functional was assumed to be continuous, for every $x \in X$ there exists a function in \mathcal{H} representing evaluation at x given by $\mathcal{R}: \delta_x \mapsto \mathcal{R}\delta_x$. In fact, it can be explicitly written:

$$\mathcal{R}\delta_x = \sum_{i \in I_2} \langle \mathcal{R}\delta_x, \phi_i \rangle_{\mathcal{H}} \phi_i = \sum_{i \in I_2} \phi_i(x) \phi_i$$

Futhermore:

$$\mathcal{R}\delta_x(x) = \sum_{i \in I_2} \phi_i(x)^2 \Rightarrow (\phi_i(x)) \in \ell^2(I_2)$$

Lemma: $\text{span}\{\mathcal{R}\delta_x\}_{x \in X}$ is dense in \mathcal{H} .

Proof: Suppose not. This is equivalent to $\text{span}\{(\phi_i(x))_{i \in I_2}\}_{x \in X}$ is not dense in $\ell^2(I_2)$.

Then exists $\bar{a} \neq 0 \in \text{span}\{(\phi_i(x))_{i \in I_2}\}_{x \in X}^\perp$ implying $\sum_{i \in I} a_i \phi_i(x) = 0 \forall x \in X$, contradicting that $\{\phi_i\}_{i \in I_2}$ is a basis.

Therefore any $f \in \mathcal{H}$ belongs to the closure of $\{\mathcal{R}\delta_x\}_{x \in X}$, and can be written in terms of the representer of the evaluation functionals. However, consider the following which motivates the mapping.

Claim: Given data $((x_1, y_1), \dots, (x_m, y_m))$, the best approximation minimum norm solution belongs to $H_{\bar{x}} := \text{span}\{\mathcal{R}\delta_{x_1}, \dots, \mathcal{R}\delta_{x_m}\}$.

Proof: Let $H = H_{\bar{x}} \oplus H_{\bar{x}}^\perp$ and for $f \in H$ factorize accordingly $f = f_{\bar{x}} + f_{\bar{x}}^\perp$. Then

$$\begin{aligned} f(x_i) &= \langle f, \mathcal{R}\delta_{x_i} \rangle_H = \langle f_{\bar{x}} + f_{\bar{x}}^\perp, \mathcal{R}\delta_{x_i} \rangle_H = \langle f_{\bar{x}}, \mathcal{R}\delta_{x_i} \rangle_H = f_{\bar{x}}(x_i) \\ \|f\|_H &= \|f_{\bar{x}} + f_{\bar{x}}^\perp\|_H = \|f_{\bar{x}}\|_H + \|f_{\bar{x}}^\perp\|_H \Rightarrow \|f_{\bar{x}}\|_H \leq \|f\|_H \end{aligned}$$

This result applied to kernels is known as the representer theorem (Mohri; Rostamizadeh; & Ameet, 2012), because it is both norm penalty and dataset loss invariant.

Because $f_{\bar{x}}$ evaluates to f at $\bar{x} = \{x_1, \dots, x_n\}$ and has smaller or equal norm, for minimum or penalized norm approximation problems at these points we only need to consider $f_{\bar{x}} \in H_{\bar{x}}$. A basis of $H_{\bar{x}}$ or linear combination $f_{\bar{x}} = \sum_{i=1}^n a_i \mathcal{R}\delta_{x_i}$ is sufficient for approximation. Because $\{\mathcal{R}\delta_{x_1}, \dots, \mathcal{R}\delta_{x_n}\}$ span the subspace by definition, but are not necessarily linearly independent, it is possible that multiple coefficient vectors $(a_i)_{i=1}^n \leftrightarrow \sum_{i=1}^n a_i \mathcal{R}\delta_{x_i}$ denote the same function. The Gram matrix $(G_{\bar{a}})_{i,j} = \left\langle \mathcal{R}\delta_{x_i}, \mathcal{R}\delta_{x_j} \right\rangle_H$ is then positive semidefinite and multiple vectors in which denote the same solution in $H_{\bar{x}}$ minimize the norm.

The problem is to find a minimum norm ($\|\bar{a}\| = \bar{a}^T G_{\bar{a}} \bar{a}$) solution to:

$$\begin{pmatrix} \mathcal{R}\delta_{x_1}(x_1) & \cdots & \mathcal{R}\delta_{x_n}(x_1) \\ \vdots & \ddots & \vdots \\ \mathcal{R}\delta_{x_1}(x_n) & \cdots & \mathcal{R}\delta_{x_n}(x_n) \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

Having obtained \bar{a} the solution is $f(x) = \sum_{i=1}^n a_i \mathcal{R}\delta_{x_i}(x)$.

Now assume we define the following symmetric positive definite function, also called a **kernel**:

$$k(x, y) = \sum_{i=1}^n \phi_i(x) \phi_i(y)$$

Rewrite the previous subspace problem using this newly defined function using the fact that $k(x_i, x_j) = \mathcal{R}\delta_{x_i}(x_j) = \left\langle \mathcal{R}\delta_{x_i}, \mathcal{R}\delta_{x_j} \right\rangle_H$. Define the matrix $K_{i,j} = k(x_i, x_j)$. The

problem is to find a minimum norm ($\|\bar{a}\| = \bar{a}^T K \bar{a}$) or solution to:

$$\begin{pmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

Having obtained $(a_i)_{i=1}^n$ solution is $f(x) = \sum_{i=1}^n a_i k(x, x_i)$.

Example

Define the function $k(x, y) = (1 + x \cdot y)^d$. Expanding, it can be written:

$$k(x, y) = 1 \cdot 1 + \sqrt{\binom{d}{1}} x \cdot \sqrt{\binom{d}{1}} y + \sqrt{\binom{d}{2}} x^2 \cdot \sqrt{\binom{d}{2}} y^2 + \dots + x^d \cdot y^d$$

Denote $\phi_0(x) = 1, \phi_1(x) = \sqrt{\binom{d}{1}}x, \phi_2(x) = \sqrt{\binom{d}{2}}x^2, \dots, \phi_d(x) = x^d$, which are clearly linearly independent. Define a Hilbert space $\mathcal{H} = \text{span}\{\phi_0, \phi_1, \phi_2, \dots, \phi_d\}$ with the norm $\|\sum_{i=0}^n a_i \phi_i\|_{\mathcal{H}}^2 = \sum_{i=0}^n a_i^2$ so that $\{\phi_i\}_{i=0}^d$ is an orthonormal basis. Then the minimum norm best approximation problem can be solved in two different ways. The regularized solutions may be obtained, for example, minimizing the Tikhonov functional as before.

Using the $n \times d$ sampling matrix: Take minimum $\|\bar{a}\|_X = \bar{a}^T \bar{a}$ norm closest approximation

$$\begin{pmatrix} \phi_0(x_1) & \dots & \phi_d(x_1) \\ \vdots & \ddots & \vdots \\ \phi_0(x_n) & \dots & \phi_d(x_n) \end{pmatrix} \begin{pmatrix} a_0 \\ \vdots \\ a_d \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

The solution is $f(x) = \sum_{i=0}^d a_i \phi_i(x)$.

Using the $n \times n$ kernel matrix. Take minimum $\|\bar{a}\|_K = \bar{a}^T K \bar{a}$ norm closest approximation

$$\begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

The solution is $f(x) = \sum_{i=1}^n \alpha_i k(x, x_i)$.

In this case the kernel is fast to compute $k(x, y) = (1 + x \cdot y)^d$ yet in most of the practical scenarios $n > d$. However, generalizing this kernel to the polynomial case $k: \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$:

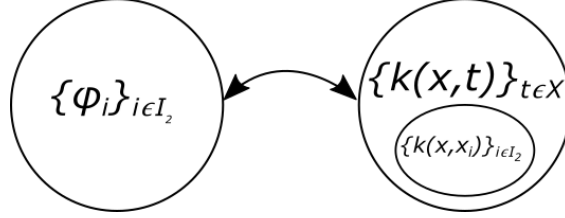
$$k(\bar{x}, \bar{y}) = (1 + \bar{x} \cdot \bar{y})^d$$

The kernel is still fast to compute, but we now have $\binom{m+d}{d}$ basis functions, rapidly increasing the computational complexity of the ordinary solution. However, if $k(\bar{x}, \bar{y})$ is computed in constant time c , the kernel matrix may still be computed in $O(cn^2)$ time and the solution then obtained in $O(n^3)$ time.

However, one should note that aside from defining the implicit possibly complicated function space, a kernel defines a norm as well for this space. For example, compared to the traditional setting of using $\{1, x, x^2, \dots, x^d\}$ as a basis of \mathcal{H} with the norm

$\|a_0 + a_1x + a_2x^2 + \dots + a_dx^d\| = \sum_{i=0}^d a_i^2$, this kernel weighted the norm by $\|a_0 + a_1x + a_2x^2 + \dots + a_dx^d\|_{\mathcal{H}} = a_0^2 + \sqrt{\binom{d}{1}}a_1^2 + \sqrt{\binom{d}{2}}a_2^2 + \dots + a_d^2$.

2.5.4 Reverse direction



We proved that if we had a Hilbert space H of functions $f: X \rightarrow \mathbb{R}$ with orthonormal basis $(\phi_i)_{i \in I}$, we can create an equivalent approximation problem by defining a symmetric positive definite function $k(x, y) = \sum_{i \in I} \phi_i(x)\phi_i(y)$ which spans the space $H_K = \text{span}\{k(x, t)\}_{t \in X}$ equipped with the norm $\langle \sum_{i=0}^d a_i k(x, t_i), \sum_{i=0}^d b_i k(x, t_j) \rangle_{H_K} = \sum_{i=1}^n \sum_{j=1}^n a_i b_j k(t_i, t_j)$. Furthermore, given a sample $\{x_1, \dots, x_n\}$ this space simplifies seeking the minimum norm solution into $H_{\bar{x}} = \text{span}\{k(x, t)\}_{t \in \{x_1, \dots, x_n\}}$.

Suppose we are given an arbitrary kernel, a positive-definite symmetric function $k(x, y)$. Which space and orthonormal basis does substituting it in the approximation problem, thus using the closure of the space $H_K = \text{span}\{k(x, t)\}_{t \in X}$ with \langle, \rangle_{H_K} , to obtain a solution correspond to? If we know a decomposition $k(x, y) = \sum_{i \in I} \phi_i(x)\phi_i(y)$ such that $\{\phi_i(x)\}_{i \in I}$ is a basis of H_K , is unknown, we may either

- Construct the closure and use H_K as the analysis space.
- Apply Gram-Schmidt in H_K to obtain a sequence of orthonormal vectors.
- Apply Mercer's theorem to obtain a decomposition with this property.

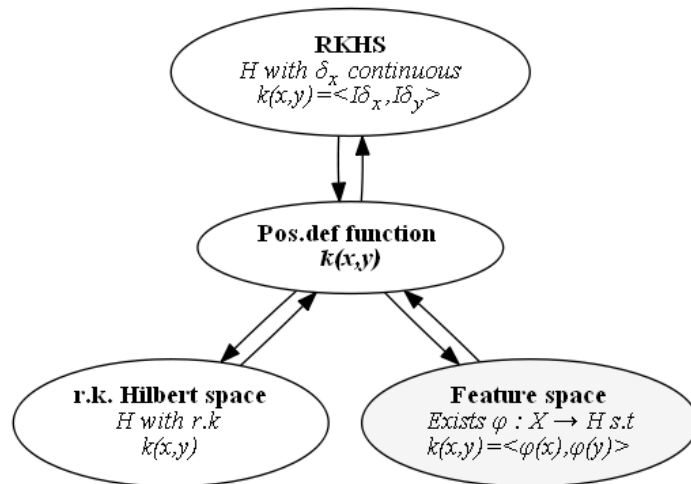
This means using a given kernel $k(x, y)$ to solve an approximation problem corresponds to implicitly using a possibly complicated, even infinite, function space. An orthonormal basis for the space can be obtained from Mercer's theorem, for example. This is advantageous compared to solving using the full basis for the function space and the sampling operator if the kernel is easy to compute and data dimensionality is high compared to the amount of data. In applications, kernels are chosen to utilize function spaces, or alternatively feature spaces, with well-known properties applicable to the task at hand and the space is taken to be implicit. The next chapter studies this based on an arbitrary kernel $k(x, y)$.

3 Reproducing Kernel Hilbert Space

3.1 Spaces associated to a kernel

To each kernel k , there corresponds an unique space of functions \mathcal{H} for which k is a reproducing kernel. Often the existence of this space is taken to be implicit, and the algorithms are expressed solely in terms of the kernel function. The solution then consists of linear combination of functions obtained by fixing one argument of the kernel: $k(x, x_i): X \rightarrow \mathbb{R}$. Yet, the space of functions still exists behind the scenes. An explicit formulation is particularly useful in our case since this allows the learning problem to be formulated as a linear operator from the kernel Hilbert space to the sample Hilbert space. Another characterization often used in machine learning is given by providing an explicit feature mapping ϕ to a feature space \mathcal{H}_0 and the kernel as the inner product in this feature space. This space may not consist of functions, but there is an unique corresponding space of functions \mathcal{H} which the learning problem in \mathcal{H}_0 corresponds to.

In the literature, there are several different characterizations of the space of functions \mathcal{H} associated to a kernel k . These characterizations are all equivalent, and are closely related to the feature space formulation. The space \mathcal{H} has several interesting properties, some following most naturally from a given characterization. These properties are used in theoretical and practical considerations, and are therefore worth illuminating. In the following we demonstrate the equivalence of 3 characterizations:



In the diagram, arrows denote implications we are going to prove. As can be seen, having shown these implications, we have equivalence between the concepts. Gray-ing of the feature space denotes a caveat; the feature space is not unique as shall be seen. One can think of the node as ‘all the feature spaces’ with the inner product equivalent to k .

(Aronszajn, Theory of Reproducing Kernels, 1948) originated much of the theory of reproducing kernels, and is the most often referred resource. An introduction to reproducing kernels with emphasis on machine learning is (Steinwart & Christmann, 2008), to which this section is mostly based on.

Definitions

In the following, let \mathcal{H} be a Hilbert space of functions $f: X \rightarrow \mathbb{R}$. Denote the inner product $\langle f, g \rangle_{\mathcal{H}}$ and the corresponding norm $\|f\|_{\mathcal{H}}$.

For a fixed $x \in X$, the map $\delta_x: \mathcal{H} \rightarrow \mathbb{R}$ defined by $\delta_x: f \mapsto f(x)$ is called the evaluation functional at x .

Definition (RKHS)

\mathcal{H} is a Reproducing Kernel Hilbert Space (RKHS) if δ_x is continuous $\forall x \in X$.

Definition (reproducing kernel)

A function $k: X \times X \rightarrow \mathbb{R}$ is called a reproducing kernel of \mathcal{H} if it satisfies:

$$\begin{aligned} \forall x \in X, \quad k(\cdot, x) &\in \mathcal{H} \\ \forall x \in X, \forall f \in \mathcal{H}, \quad \langle f, k(\cdot, x) \rangle_{\mathcal{H}} &= f(x) \end{aligned}$$

Definition (Feature space)

The function $k: X \times X \rightarrow \mathbb{R}$ is a kernel if there exists a real Hilbert space \mathcal{F} and a map $\phi: X \rightarrow \mathcal{F}$ such that

$$\forall x, y \in X, k(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{F}}$$

The map ϕ is called a feature map and \mathcal{F} a feature space. Note that in the definition, we do not require \mathcal{F} to be a space of functions; hence a different symbol.

Definition (positive definite function)

A symmetric function $h: X \times X \rightarrow \mathbb{R}$ is positive definite if

$$\forall n \geq 1, \bar{a} \in \mathbb{R}^n, \forall \bar{x} \in X^n, \sum_{i=1}^n \sum_{j=1}^n a_i a_j h(x_i, x_j) \geq 0$$

Equivalently, we may require the matrix $H_{ij} := h(x_i, x_j)$ to be positive semidefinite for all $n \in \mathbb{N}, x_1, \dots, x_n$.

3.2 Equivalence of spaces

Theorem

Given a symmetric positive definite function $k: X \times X \rightarrow \mathbb{R}$, define the following space of functions (Steinwart & Christmann, 2008):

$$\mathcal{H}_{pre} := \left\{ \sum_{i=1}^n \alpha_i k(\cdot, x_i) : \alpha_1, \dots, \alpha_n \in \mathbb{R}, x_1, \dots, x_n \in X \right\}$$

Given $f := \sum_{i=1}^n \alpha_i k(\cdot, x_i)$ and $g := \sum_{j=1}^m \beta_j k(\cdot, x'_j)$, define:

$$\langle f, g \rangle := \sum_{i=1}^n \sum_{j=1}^m \alpha_i \beta_j k(x_i, x'_j)$$

Then \mathcal{H}_{pre} is an inner product space. Denote \mathcal{H} the completion, which is then a Hilbert space with a reproducing kernel k .

Futhermore the reproducing kernel Hilbert space is unique. If k is the reproducing kernel of \mathcal{F} , then an isometric isomorphism exists between \mathcal{F} and \mathcal{H} .

Proof

$\langle \cdot, \cdot \rangle$ is a well-defined inner product on \mathcal{H}_{pre} :

- Well-defined: $\langle f, g \rangle = \sum_{j=1}^m \beta_j f(x'_j) = \sum_{i=1}^n \alpha_i g(x_i)$, implying independent of representation of f and g .
- Bilinear, symmetric and positive: follows directly with k assumed symmetric positive definite.
- Separates points: Schwarz inequality implies $|\langle f, g \rangle| = \langle f, f \rangle \langle g, g \rangle \forall f, g \in \mathcal{H}_{pre}$. Then $\|f\| = 0$ implies $|f(x)| = |\sum_{i=1}^n \alpha_i k(x, x_i)|^2 = |\langle f, k(\cdot, x) \rangle|^2 \leq \langle f, f \rangle \langle k(\cdot, x), k(\cdot, x) \rangle = \|f\|^2 k(x, x) = 0$.

Then k is the reproducing kernel of \mathcal{H}_{pre} :

- $k(\cdot, x) \in \mathcal{H}_{pre} \forall x \in X$

- $\langle f, k(\cdot, x) \rangle = f(x) \forall f \in \mathcal{H}_{pre}$

\mathcal{H}_{pre} is a well-defined inner product space with the reproducing kernel k , but it may not be a Hilbert space, i.e. complete.

Completion is defined as an extension of a space such that the new space is complete and contains a dense subspace isometric with the original space. There exists a standard completion of any inner product space $\overline{\mathcal{H}_{pre}} = \mathcal{H}$ (Kreyszig, 1989, ss. 41,69,139). Given a space \mathcal{H}_{pre} with inner product $\langle \cdot, \cdot \rangle$ and induced norm $\|\cdot\|$, completion is defined:

$$\mathcal{H} = \{[(f_n)] : (f_n) \text{ is Cauchy in } \mathcal{H}_{pre}\}$$

Where $[(f_n)]$ is the equivalence class of (f_n) defined relative to the equivalence $(f_n) \sim (g_n)$ if

$$\lim_{n \rightarrow \infty} \|f_n - g_n\| = 0$$

To make the space a vector space, addition and scalar multiplication are defined:

$$\begin{aligned} [(f_n)] + [(g_n)] &= [(f_n + g_n)] \\ \alpha[(f_n)] &= [(\alpha f_n)] \end{aligned}$$

For any representatives $(f_n) \in [(f_n)]$ and $(g_n) \in [(g_n)]$ an inner product on \mathcal{H} is defined:

$$\langle [(f_n)], [(g_n)] \rangle_{\mathcal{H}} = \lim_{n \rightarrow \infty} \langle f_n, g_n \rangle$$

It can be shown that this construction satisfies the properties required from completion. The inner product and vector space operations are well defined and independent of the representers. The mapping $T: \mathcal{H}_{pre} \rightarrow \mathcal{H}$ defined by $f \mapsto [(f)]$ is an isometric mapping onto a dense subspace of \mathcal{H} . Furthermore, \mathcal{H} is complete and unique except for isometries.

However, the space \mathcal{H} does not consist of functions. When a space of functions is completed such that the completion consists of functions as well and preserves certain continuity properties, we speak of **functional completion** (Aronszajn, Theory of Reproducing Kernels, 1948). For each Cauchy sequence (f_n) in \mathcal{H}_{pre} :

$$\begin{aligned} |f_m(x) - f_n(x)| &= |\langle f_m, k(\cdot, x) \rangle - \langle f_n, k(\cdot, x) \rangle| = |\langle f_m - f_n, k(\cdot, x) \rangle| \\ &\leq \|f_m - f_n\| k(x, x)^{\frac{1}{2}} \end{aligned}$$

Since $(f_n(x))$ is Cauchy in \mathbb{R} which is complete, the sequence converges. Substituting two equivalent sequences in the inequality we see $(f_n) \sim (g_n) \Rightarrow |f_n(x) - g_n(x)| \rightarrow 0$. Therefore it is possible to associate to each equivalence class $[(f_n)]$ the function:

$$f(x) = \lim_{n \rightarrow \infty} f_n(x)$$

Consider the sequence $(k(\cdot, x))$ associated to $k(\cdot, x)$. Existence of the reproducing kernel is easily verified:

$$f(x) = \lim_{n \rightarrow \infty} f_n(x) = \lim_{n \rightarrow \infty} \langle f_n, k(\cdot, x) \rangle = \langle [(f_n)], [(k(\cdot, x))] \rangle_{\mathcal{F}}$$

Given an arbitrary \mathcal{F} with a reproducing kernel k , first property $k(\cdot, x) \in \mathcal{F} \forall x$ implies $\mathcal{H}_{pre} \subset \mathcal{F}$. The reproducing property implies $\forall f \in (\mathcal{H}_{pre})^\perp: f(x) = \langle f, k(\cdot, x) \rangle = 0 \Rightarrow (\mathcal{H}_{pre})^\perp = \{0\}$ which implies \mathcal{H}_{pre} is dense in \mathcal{F} . Because the completion is unique up to isometric isomorphism (Kreyszig, 1989), we have $\mathcal{H} = \mathcal{F}$.

Finally, every space has a unique reproducing kernel, since if k_1 and k_2 are reproducing kernels of \mathcal{H} :

$$\begin{aligned} \|k_1(\cdot, x) - k_2(\cdot, x)\|^2 &= \langle k_1(\cdot, x) - k_2(\cdot, x), k_1(\cdot, x) - k_2(\cdot, x) \rangle \\ &= k_1(x, x) - k_2(x, x) - (k_1(x, x) - k_2(x, x)) = 0 \end{aligned}$$

It follows $\forall x: k_1(y, x) = k_2(y, x) \Rightarrow k_1 = k_2$.

This shows that to every symmetric positive definite function $k(x, y)$ corresponds a unique Hilbert space of functions for which k is a reproducing kernel. For a space with reproducing kernel k , the kernel is unique and symmetric positive semidefinite.

Theorem

\mathcal{H} is an RKHS if and only if it has a reproducing kernel. Furthermore, the reproducing kernel is given by (Steinwart & Christmann, 2008):

$$k(x, y) = \langle \mathcal{R}\delta_y, \mathcal{R}\delta_x \rangle$$

Where $\mathcal{R}: \mathcal{H}' \rightarrow \mathcal{H}$ maps a bounded functional to the representer in \mathcal{H} , given by Riesz's theorem (Kreyszig, 1989, s. 188).

Proof

k is the reproducing kernel of \mathcal{H} :

- $k(x, y) = \langle \mathcal{R}\delta_y, \mathcal{R}\delta_x \rangle = \delta_x(\mathcal{R}\delta_y) = \mathcal{R}\delta_y(x) \Rightarrow k(\cdot, y) = \mathcal{R}\delta_y \in \mathcal{H} \forall y$
- $f(x) = \delta_x(f) = \langle f, \mathcal{R}\delta_x \rangle = \langle f, k(\cdot, x) \rangle$

If \mathcal{H} has a reproducing kernel k , δ_x is continuous:

$$|\delta_x(f)| = |f(x)| = |\langle f, k(\cdot, x) \rangle| \leq \|f\| \|k(\cdot, x)\|^{\frac{1}{2}}$$

RKHS formulation is often used to highlight an important convergence property (Aronszajn, Theory of Reproducing Kernels, 1948). Let (f_n) be a sequence converging to f : $\|f_n - f\| \rightarrow 0$ as $n \rightarrow \infty$. Then continuity of δ_x means norm convergence implies pointwise convergence:

$$\lim_{n \rightarrow \infty} f_n(x) = \lim_{n \rightarrow \infty} \delta_x(f_n) = \delta_x(f) = f(x)$$

Theorem

$k: X \times X \rightarrow \mathbb{R}$ is a symmetric positive definite function if and only if it is a kernel.

Proof

Let \mathcal{H} be a Hilbert space for which the symmetric positive definite function k is a reproducing kernel. Define the feature mapping $\Phi: X \rightarrow \mathcal{H}$, also called the *canonical feature map* (Steinwart & Christmann, 2008):

$$\Phi(x) = k(\cdot, x)$$

Then k realizes the inner product on \mathcal{H} :

$$\langle \Phi(x), \Phi(y) \rangle_{\mathcal{H}} = \langle \Phi(y), \Phi(x) \rangle_{\mathcal{H}} = \langle k(\cdot, y), k(\cdot, x) \rangle_{\mathcal{H}} = k(x, y)$$

Suppose $k(x, y) = \langle \Phi(x), \Phi(y) \rangle_{\mathcal{H}}$. Then k is symmetric positive definite:

$$\sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(x_i, x_j) = \left\langle \sum_{i=1}^n \alpha_i \Phi(x_i), \sum_{j=1}^n \alpha_j \Phi(x_j) \right\rangle_{\mathcal{H}} \geq 0$$

Note that the theorem implies existence of a feature map, not uniqueness. There are several feature spaces for which kernel realizes an inner product. For a trivial example, define the feature maps $\phi_1: x \mapsto x \in \mathbb{R}$ and $\phi_2: x \mapsto \left(\frac{x}{\sqrt{2}}, \frac{x}{\sqrt{2}}\right) \in \mathbb{R}^2$. Then

$k(x, y) = xy = \langle \phi_1(x), \phi_1(y) \rangle_{\mathbb{R}} = \langle \phi_2(x), \phi_2(y) \rangle_{\mathbb{R}^2}$. However as seen above, the function space for which k is a reproducing kernel is unique, there exists a feature map which realizes it as a feature space. Furthermore, the following fact implies that there exists a metric surjection from any feature space into it, which means it can be seen as the *smallest* feature space in this sense.

Theorem (Steinwart & Christmann, 2008)

Let $k: X \times X \rightarrow \mathbb{R}$ be a kernel with feature space \mathcal{H}_0 and feature map $\Phi_0: X \rightarrow \mathcal{H}_0$. Then the RKHS may also be expressed through the arbitrary feature space:

$$\mathcal{H} = \{f: X \rightarrow \mathbb{R} : \exists w \in \mathcal{H}_0 \ f(x) = \langle w, \Phi(x) \rangle_{\mathcal{H}_0}\}$$

equipped with the norm:

$$\|f\|_{\mathcal{H}} := \inf\{\|w\|_{\mathcal{H}_0} : w \in \mathcal{H}_0 \ f(x) = \langle w, \Phi(x) \rangle_{\mathcal{H}_0}\}$$

Moreover, a metric surjection $V: \mathcal{H}_0 \rightarrow \mathcal{H}$ may be defined by

$$(Vw)(x) = \langle w, \Phi(x) \rangle_{\mathcal{H}_0} \ \forall w \in \mathcal{H}_0$$

This theorem is particularly useful for considering the effect on the RKHS of modifications of symmetric positive definite functions by compositions through operations on the feature space, compared to the complicated proofs in (Aronszajn, Theory of Reproducing Kernels, 1948).

3.3 Kernel and space modifications

Consider kernels K_i with RKHSs \mathcal{H}_i . If we have a composition K^* consisting of K_i which can be readily verified to be positive definite, we know that there exists an RKHS \mathcal{H}^* associated to it. Likewise, if an operation such as restriction modifies kernel K with RKHS \mathcal{H} such that K remains symmetric positive semidefinite there is a RKHS $\mathcal{O}(\mathcal{H})$ associated to the modified kernel $\mathcal{O}(K)$. Can these RKHS be expressed in terms of the original \mathcal{H} ? (Aronszajn, 1948) systematically developed the following results, with the proofs expressed much more simply using the theorem in the previous section:

Sum of Kernels

If K_1, K_2 are reproducing kernels of $\mathcal{H}_1, \mathcal{H}_2$, then $K(x, y) = K_1(x, y) + K_2(x, y)$ is the reproducing kernel of $\mathcal{H} = \{f_1 + f_2 \mid f_1 \in \mathcal{H}_1, f_2 \in \mathcal{H}_2\}$ with the norm defined as $\|f\|_{\mathcal{H}}^2 = \min_{f_1 \in \mathcal{H}_1, f_2 \in \mathcal{H}_2} \|f_1\|_{\mathcal{H}_1}^2 + \|f_2\|_{\mathcal{H}_2}^2$.

Proof

Let $K_1(x, y) = \langle \phi_1(x), \phi_1(y) \rangle_{\mathcal{H}_1}, K_2(x, y) = \langle \phi_2(x), \phi_2(y) \rangle_{\mathcal{H}_2}$. Define the map $\phi(x) = (\phi_1(x), \phi_2(x))$. Then $\phi: X \rightarrow \mathcal{H}_1 \oplus \mathcal{H}_2$ is a feature map with kernel $K(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}} = \langle \phi_1(x), \phi_1(y) \rangle_{\mathcal{H}_1} + \langle \phi_2(x), \phi_2(y) \rangle_{\mathcal{H}_2} = K_1(x, y) + K_2(x, y)$. Since $\langle w, \phi(y) \rangle_{\mathcal{H}} = \langle w_1, \phi_1(y) \rangle_{\mathcal{H}_1} + \langle w_2, \phi_2(y) \rangle_{\mathcal{H}_2}$ substituting this to the theorem gives the RKHS result.

Product of Kernels

If K_1, K_2 are reproducing kernels of $\mathcal{H}_1, \mathcal{H}_2$, then $K(x_1, x_2, y_1, y_2) = K_1(x_1, y_1)K_2(x_2, y_2)$ is the reproducing kernel of $\mathcal{H} = \mathcal{H}_1 \otimes \mathcal{H}_2$.

Proof

Let $K_1(x, y) = \langle \phi_1(x), \phi_1(y) \rangle_{\mathcal{H}_1}, K_2(x, y) = \langle \phi_2(x), \phi_2(y) \rangle_{\mathcal{H}_2}$. Define the map $\phi(x) = \phi_1(x) \otimes \phi_2(x)$. Then $\phi: X \rightarrow \mathcal{H}_1 \otimes \mathcal{H}_2$ is a feature map with kernel $K(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}} = \langle \phi_1(x) \otimes \phi_2(x), \phi_1(y) \otimes \phi_2(y) \rangle_{\mathcal{H}} = \langle \phi_1(x), \phi_1(y) \rangle_{\mathcal{H}_1} \langle \phi_2(x), \phi_2(y) \rangle_{\mathcal{H}_2} = K_1(x, y)K_2(x, y)$. Substituting $\langle w, \phi(y) \rangle_{\mathcal{H}} = \langle w_1, \phi_1(y) \rangle_{\mathcal{H}_1} \langle w_2, \phi_2(y) \rangle_{\mathcal{H}_2}$ to the theorem gives the RKHS result.

Taylor series of a Kernel (Steinwart & Christmann, 2008)

Denote $\dot{B}_{\mathbb{R}}$ and $\dot{B}_{\mathbb{R}^d}$ the unit balls in \mathbb{R} and \mathbb{R}^d and let $r \in [0, \infty)$. Assume f is analytic with Taylor series $f(x) = \sum_{i=0}^{\infty} a_i z^i$ where $z \in r\dot{B}_{\mathbb{R}}$ and $a_i \geq 0$. Then

$$k(x, y) = f(\langle x, y \rangle) = \sum_{i=0}^{\infty} a_i \langle x, y \rangle^i \text{ for } x, y \in \sqrt{r}\dot{B}_{\mathbb{R}^d}$$

Is a kernel, called **taylor type** kernel.

Restriction of Kernels

If K is the reproducing kernel of \mathcal{H} of functions $f: X \rightarrow \mathbb{R}$, then $K|_{X'}$, defined K restricted to $X' \subset X$ is the reproducing kernel of \mathcal{H}_0 of functions $f|_{X'}$ where $f \in \mathcal{H}$ with the norm $\|f\|_{\mathcal{H}_0}^2 = \min_{g \in \mathcal{H}: g|_{X'} = f} \|g\|_{\mathcal{H}}^2$.

Proof

Let $K(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}}$. Define the map $\phi_0(x) = \phi(Px)$ where $P: X' \rightarrow X$ is inclusion. Then $\phi_0: X' \rightarrow \mathcal{H}$ is a feature map with kernel $K_{|X'}(x, y) = \langle \phi_0(x), \phi_0(y) \rangle_{\mathcal{H}_0} = \langle \phi(Px), \phi(Py) \rangle_{\mathcal{H}} = K(Px, Py)$. Substituting $\langle w_0, \phi_0(y) \rangle_{\mathcal{H}_0} = \langle w, \phi(Py) \rangle_{\mathcal{H}}$ to the theorem gives the RKHS result.

These simple results can be used to construct many commonly used kernels on \mathbb{R}^d , for example the following are often used:

- Linear kernel: $k(x, y) = \langle x, y \rangle$
- Polynomial kernel: $k(x, y) = (\langle x, y \rangle + c)^m$ where $c \geq 0$.
- Exponential kernel: $k(x, y) = \exp(2\sigma x \cdot y)$
- Gaussian RBF kernel: $k(x, y) = \exp(-\sigma \|x - x'\|^2)$

3.4 Prior knowledge in Kernels

Multiple methods for incorporating prior knowledge into the kernel learning process have been considered in the literature. Prior knowledge is classified into two main categories in (Lauer & Bloch, 2007): invariance in the input space and knowledge on the data. By invariance they refer to the constraint that the output should not change under transformations or permutations of samples; for example transition or rotation should not change the letter in digit recognition. They define knowledge on the data as manipulating the data set: one may incorporate additional knowledge on unlabelled data such as training set imbalance or add quality relative weighting on the labelled sample.

Formally, invariance in the input space can be expressed as a constraint in terms of a parametrized transformation $T_\theta: x \mapsto T_\theta x$ of the input pattern (Lauer & Bloch, 2007):

$$f(x) = f(T_\theta x) \quad \forall x, \theta$$

Often only local invariance is enforced through the tangent at $\theta = 0$:

$$\frac{\delta}{\delta \theta} \big|_{\theta=0} f(T_\theta x)$$

Some methods constrain the domain directly without specifying a transformation:

$$f(x) = y_{\mathcal{P}} \quad \forall x \in \mathcal{P} \text{ where } \mathcal{P} \subseteq X \text{ is constrained}$$

These constraints may be enforced in different places of the learning setting. (Schölkopf & Smola, 2001) classify the methods into three types based on the component they affect:

1. Sample methods modify data by generating or modifying samples
2. Kernel methods modify the kernel function
3. Optimization methods modify the minimization process or problem formulation

A popular sample method is virtual examples (Decoste & Bernhard, 2002). This method simply generates new samples from the training data which satisfy the invariance and have the same label. For a singular non-parametric transformation new example is added for each vector in the training set $(x_i, y_i) \mapsto (Tx_i, y_i)$, $i = 1 \dots N$. If the transformation is continuous, the transformation may be discretized and samples added for $T_\theta x_i, y_i, \theta = 1, \dots M$. This method is used in neural networks as hints and in SVMs as virtual SVs. For enforcing symmetry in the training data it is particularly simple, because we simply need to add

$$((v, v'), y_i) \mapsto ((v', v), y_i)$$

It can be shown that for permutation invariant kernels and non-zero regularization parameter it results in enforcing symmetry in the solution as well.

The method used in our study is to modify the kernel to incorporate prior knowledge. One example of such knowledge is invariance of kernel k under given transition T :

$$k(x, z) = k(Tx, z)$$

The review (Lauer & Bloch, 2007) mentions at least four methods which utilize the same technique: jittering kernels, tangent distance (TD), Haar-integration (HI) and Kernels between sets. In our case, the transition is not arbitrary like in these methods and the simplicity allows additional considerations. Because these methods are useful illustrations for their generality, they are briefly reviewed.

3.5 Examples of prior knowledge in Kernels

Jittering kernels (Decoste & Bernhard, 2002)

The jittered variations (translated, rotated, etc.) of samples are incorporated to the kernel itself. The new kernel $k^J(x_i, x_j)$ is defined in terms of the old kernel as follows

- Select a jittered variation $x_q \in \mathcal{L}(x_i)$ of a sample x_i closest to x_j by minimizing the distance $\|x_q - x_j\|_{\mathcal{H}} = \sqrt{k_{qq} - 2k_{qj} + k_{jj}}$ in the feature space \mathcal{H} .
- Define a new kernel by $k_{ij}^J = k_{qj}$

Tangent distance (TD) (Haasdonk & Keysers, 2002)

Let P_x be the set of all samples generated by transformation T from a sample x . For sample x_1 and x_2 tangent distance uses approximation techniques to define the kernel through the distance between P_{x_1} and P_{x_2} . Often linear approximation of the transformation T is used, which may be extended by using the first terms in the Taylor expansion of T . Express Tx as the combination of n_L local and simple parametrized linear transformations L_{α_k} , $\alpha_k \in [\alpha_k^{\min}, \alpha_k^{\max}]$ based on tangent vectors ℓ_{α_k} :

$$Tx \approx x + \sum_{k=1}^{n_L} \alpha_k L_{\alpha_k}(x)$$

TD is then defined as the minimal Euclidean distance between all approximated transformations:

$$\rho(P_x, P_z) = \left(x - z + \sum_{k=1}^{n_L} \left(\alpha_k \ell_{\alpha_k}(x) - \beta_k \ell_{\beta_k}(z) \right) \right)^2, \alpha_k \in [\alpha_k^{\min}, \alpha_k^{\max}], \beta_k \in [\beta_k^{\min}, \beta_k^{\max}]$$

The input space tangent distance is then combined with a kernel which uses the Euclidean distance in the input space, for example the Gaussian kernel $k(x, y) = \exp\left(-\frac{\rho(P_x, P_z)^2}{2\sigma^2}\right)$.

Haar-integration (Schulz-Mirbach, 1994)

For a standard kernel k_0 and a transformation group \mathcal{T} , an average of the standard kernel output is computed over all pairwise combinations of the transformed examples:

$$k(x, z) = \int_{\mathcal{T}} \int_{\mathcal{T}} k_0(Tx, T'z) dT dT'$$

This is equivalent to computing the inner product between averages $\overline{\Phi(x)}$ and $\overline{\Phi(z)}$ of the transformed examples:

$$k(x, z) = \int_{\mathcal{T}} \int_{\mathcal{T}} \langle \Phi(Tx), \Phi(T'z) \rangle dT dT' = \left\langle \int_{\mathcal{T}} \Phi(Tx) dT, \int_{\mathcal{T}} \Phi(T'z) dT' \right\rangle$$

Kernels between sets (Wolf & Shashua, 2003)

Kernels between sets aim to define kernels invariant to a certain kind of permutation of the input. For example, if one considers kernels over sets of vectors, and represents these as matrices $A = [\Phi(a_1), \dots, \Phi(a_n)]$, $B = [\Phi(b_1), \dots, \Phi(b_n)]$, one may define a kernel by the principal angles θ_i in the feature space:

$$k(A, B) = \prod_{i=1}^n \cos(\theta_i)$$

This kernel does not modify an existing kernel, but it demonstrates a kernel with an invariance property with respect to the permutations of the columns of the matrices.

4 Theory of Linear Inverse Problems

In the regularized least squares formulation, the sampling operator $S_{\bar{x}}: \mathcal{H} \rightarrow \ell^2(I_1)$ is a map from a normed space $\|\cdot\|_{\mathcal{H}}$ to a normed space $\|\cdot\|_2$. In practise the sample space is often \mathbb{R}^n , which is finite dimensional. We saw that the inverse may not exist, so an optimal solution was defined through these norms:

$$f^* = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \|S_{\bar{x}}f - \bar{y}\|_2 \text{ of minimum/penalized } \|f\|_{\mathcal{H}}$$

Unlike most other machine learning algorithms (Schölkopf & Smola, 2001), using optimization techniques is not necessary because this formulation allows derivation of the solution analytically. Analytical treatment allows further insights into the algorithm. In the general, consider an abstract linear map $T: \mathcal{X} \rightarrow \mathcal{Y}$ between Hilbert spaces \mathcal{X} and \mathcal{Y} . Under certain conditions, it is then possible to define an “inverse” operator where T is not bijective. These operators are needed when kernel approximation is later considered from learning theoretic perspective.

4.1 Linear operators

Linear operator is a map $T: \mathcal{X} \rightarrow \mathcal{Y}$ between vector spaces over the same field K such that $\forall x, y \in X, \alpha \in K$:

$$T(x + y) = Tx + Ty$$

$$T(\alpha x) = \alpha Tx$$

Bounded operator is a linear operator $T: \mathcal{X} \rightarrow \mathcal{Y}$ between normed spaces such that

$$\exists c \in \mathbb{R} : \forall x \in \mathcal{X} \|Tx\|_{\mathcal{Y}} \leq c\|x\|_{\mathcal{X}}$$

Compact operator (Kreyszig, 1989) is a linear operator $T: \mathcal{X} \rightarrow \mathcal{Y}$ between normed spaces such that

$$M \subset \mathcal{X} \text{ bounded} \Rightarrow \overline{T(M)} \text{ compact}$$

A foundational result of linear algebra is that any linear operator $T: \mathcal{X} \rightarrow \mathcal{Y}$ between **finite-dimensional** vector spaces $\dim(\mathcal{X}) = m, \dim(\mathcal{Y}) = n$ with respect given bases may be represented by an $n \times m$ **matrix**:

$$T = \begin{pmatrix} t_{11} & \cdots & t_{1m} \\ \vdots & \ddots & \vdots \\ t_{n1} & \cdots & t_{nm} \end{pmatrix}$$

Combining several straightforward results in (Kreyszig, 1989), we obtain the following hierarchy:

$$\begin{aligned} \text{Finite-dimensional operator (matrix)} &\subset \text{Compact operator} \\ &\subset \text{Bounded operator} \subset \text{Linear operator} \end{aligned}$$

This hierarchy is useful because the increasingly specified cases allow increasingly powerful results.

Finally, the following linear operator is very useful for considering the spaces associated to a linear operator.

Definition (Kreyszig, 1989)

Let $T: \mathcal{X} \rightarrow \mathcal{Y}$ be a bounded Linear operator, where \mathcal{X} and \mathcal{Y} are Hilbert spaces. Then the Hilbert-**adjoint operator** T^* of T is the operator $T^*: \mathcal{Y} \rightarrow \mathcal{X}$ such that for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$:

$$\langle Tx, y \rangle_{\mathcal{Y}} = \langle x, T^*y \rangle_{\mathcal{X}}$$

4.2 Fundamental subspaces of a linear operator

Define the following subspaces, known as range and nullspace, for a linear operator $A: \mathcal{X} \rightarrow \mathcal{Y}$:

$$\begin{aligned} \mathcal{R}(A) &= \{y : y = Ax, x \in \mathcal{X}\} \\ \mathcal{N}(A) &= \{x : Ax = 0, x \in \mathcal{X}\} \end{aligned}$$

Then there are four subspaces associated to the operator $A: \mathcal{X} \rightarrow \mathcal{Y}$ and its adjoint $A^*: \mathcal{Y} \rightarrow \mathcal{X}$:

$$\begin{aligned} \mathcal{R}(A) &\subset \mathcal{Y} \\ \mathcal{N}(A) &\subset \mathcal{X} \\ \mathcal{R}(A^*) &\subset \mathcal{X} \\ \mathcal{N}(A^*) &\subset \mathcal{Y} \end{aligned}$$

With the following property:

$$\begin{aligned} \mathcal{R}(A)^\perp &= \mathcal{N}(A^*) \\ \mathcal{R}(A^*)^\perp &= \mathcal{N}(A) \end{aligned}$$

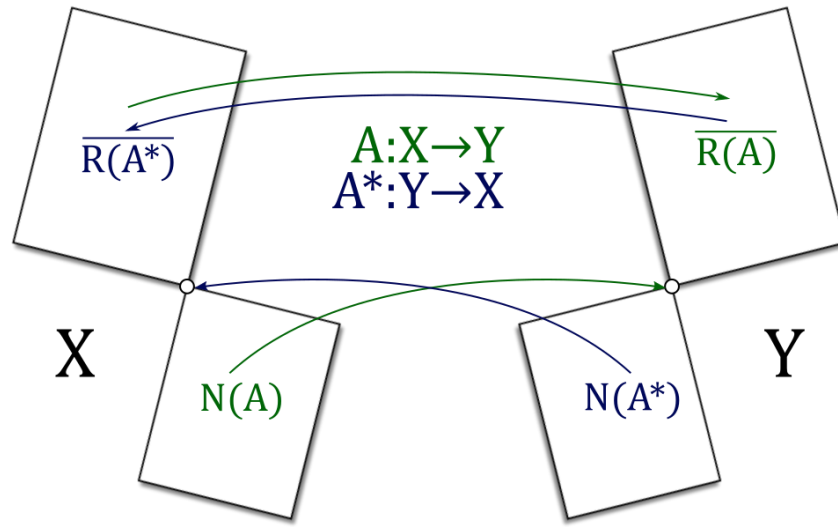
Since $\mathcal{N}(A)$ and $\mathcal{N}(A^*)$ are closed we have the factorization $\mathcal{X} = \mathcal{N}(A) \oplus \mathcal{N}(A)^\perp$ and $\mathcal{Y} = \mathcal{N}(A^*) \oplus \mathcal{N}(A^*)^\perp$. Furthermore, if $\mathcal{R}(A)$ is closed $(\mathcal{R}(A)^\perp)^\perp = \mathcal{R}(A)$, otherwise $(\mathcal{R}(A)^\perp)^\perp = \overline{\mathcal{R}(A)}$ (Kreyszig, 1989).

Therefore we get the following factorization, where closure can be omitted if the range is closed:

$$\mathcal{X} = \overline{\mathcal{R}(A^*)} \oplus \mathcal{N}(A)$$

$$\mathcal{Y} = \overline{\mathcal{R}(A)} \oplus \mathcal{N}(A^*)$$

Any vector $\bar{x} \in \mathcal{X}$ can in this case be factored into $\overline{\mathcal{R}(A^*)}$ and $\mathcal{N}(A)$ component and the action of $A\bar{x}$ expressed as a bijection $A|_{\overline{\mathcal{R}(A^*)}}$ and null mapping $A|_{\mathcal{N}(A)}$ (Moon, 1999):



$$\bar{x} = \bar{x}_r + \bar{x}_n$$

$$A\bar{x} = A\bar{x}_r + A\bar{x}_n \text{ with } A\bar{x}_r \in \mathcal{R}(A), A\bar{x}_n = 0$$

Subspaces associated to a matrix.

For a matrix $A: \mathbb{R}^n \rightarrow \mathbb{R}^m$ since $\dim(\mathcal{R}(A)) = \dim(\mathcal{R}(A^*))$ this means we have the following rank results:

$$\dim(\mathcal{R}(A)) = r$$

$$\dim(\mathcal{R}(A^*)) = r$$

$$\dim(\mathcal{N}(A)) = n - r$$

$$\dim(\mathcal{N}(A^*)) = m - r$$

Where $\mathcal{R}(A)$ is called column space, $\mathcal{N}(A)$ the nullspace and $\mathcal{R}(A^*)$ the row space. The spaces are equipped with the inner product $\langle \bar{x}, \bar{y} \rangle = \bar{x}^T \bar{y}$ implying for the adjoint $A^* = A^T$ and the terminology follows. (Kreyszig, 1989).

4.3 Spectral Theorems

The following factorizations are defined for linear operators $T: \mathcal{H} \rightarrow \mathcal{H}$ which are self-adjoint $T^* = T$, widely referred to as *spectral theorems*. They are presented in the order of abstraction.

Spectral Theorem for Bounded Operators (Kreyszig, 1989)

Let \mathcal{H} be a complex Hilbert space and $T: \mathcal{H} \rightarrow \mathcal{H}$ be a bounded self-adjoint operator. Then T has the spectral representation

$$T = \int_m^M \lambda dE_\lambda$$

where $\{E_\lambda\}$ is the spectral family associated with T . The integral is to be understood in the sense of uniform operator convergence.

This theorem is useful to illustrate the most general case, but is not used in this thesis. For details, see the details and extensive derivation in (Kreyszig, 1989, ss. 459-522).

Spectral Theorem for Compact Operators (Kreyszig, 1989)

Let \mathcal{H} be a Hilbert space and $T: \mathcal{H} \rightarrow \mathcal{H}$ be a compact self-adjoint operator. Then there is a countable orthonormal set $\{e_j\}_{j \in J}$ in \mathcal{H} and $\{\lambda_j\}_{j \in J}$ $|\lambda_1| \geq |\lambda_2| \geq \dots > 0$ converging to zero such that

$$Tf = \sum_{j \in J} \lambda_j \langle f, e_j \rangle_{\mathcal{H}} e_j$$

Diagonalization of Symmetric Matrices (Moon, 1999)

Every $m \times m$ real matrix $A: \mathbb{R}^m \rightarrow \mathbb{R}^m$ such that $A^T = A$ is called symmetric (self-adjoint in \mathbb{R}^m) and can be factored

$$A = U \Sigma U^T$$

where $U \in \mathbb{R}^{m \times m}$ is orthogonal and $\Sigma \in \mathbb{R}^{m \times m}$ has the form $\Sigma = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$. U contains the eigenvectors of A as columns and (λ_i) the corresponding eigenvalues.

4.4 Singular Value Decompositions

4.4.1 SVD for compact operators

Compact operators are an important special case, since specific integral operators are compact, and they are used in connection with kernels.

Theorem: Compact integral operators (Engl; Hanke; & Neubauer, 1996)

If $\Omega \subseteq \mathbb{R}^n$ is compact and Jordan measurable and the kernel k is continuous on $\{(s, t) \in \Omega \times \Omega \mid s \neq t\}$ and for all $s \neq t \in \Omega$, $|k(s, t)| \leq \frac{M}{|s-t|^{n-\epsilon}}$ with $M > 0$, $\epsilon > 0$, then

$$K: \mathcal{L}^2(\Omega) \rightarrow \mathcal{L}^2(\Omega)$$

$$x \mapsto (Kx)(s) := \int_{\Omega} k(s, t)x(t)dt$$

is compact.

The following two results are analogous between compact operators and finite-dimensional operators.

Theorem: Singular values of a compact operator (Engl; Hanke; & Neubauer, 1996)

With a singular system $(\{\sigma_i\}, \{v_i\}, \{u_i\})$ of orthonormal eigenvectors v_i of K^*K , orthonormal eigenvectors u_i of KK^* and nonzero eigenvalues σ_i^2 of K^*K/KK^* ordered by multiplicity, a compact operator K may be expressed:

$$\begin{aligned} Kv_i &= \sigma_i u_i \\ K^* u_i &= \sigma_i v_i \\ Kx &= \sum_{i=1}^{\infty} \sigma_i \langle x, v_i \rangle u_i \\ K^* x &= \sum_{i=1}^{\infty} \sigma_i \langle x, u_i \rangle v_i \\ K^* Kx &= \sum_{i=1}^{\infty} \sigma_i^2 \langle x, v_i \rangle v_i \\ KK^* x &= \sum_{i=1}^{\infty} \sigma_i^2 \langle x, u_i \rangle u_i \end{aligned}$$

where these infinite series converge, they are called the singular value expansions.

K has finitely many singular values iff it has a finite-dimensional range. For an integral operator K with square integrable kernel k , this is the case iff k is degenerate (Engl;Hanke;& Neubauer, 1996):

$$k(s, t) = \sum_{i=1}^n \phi_i(s) \psi_i(t)$$

If there are infinitely many singular values, they accumulate at 0 (Kreyszig, 1989):

$$\lim_{i \rightarrow \infty} \sigma_i = 0$$

As a simple application of Banach's open mapping theorem in (Engl;Hanke;& Neubauer, 1996), the range of K is closed iff it is finite-dimensional. Therefore $\mathcal{R}(K)$ of non-degenerate K is non-closed.

4.4.2 SVD for matrices

Every matrix $A \in \mathbb{R}^{m \times n}$ can be factored as

$$A = U \Sigma V^T$$

where $U \in \mathbb{R}^{m \times m}, V \in \mathbb{R}^{n \times n}$ are orthogonal and $\Sigma \in \mathbb{R}^{m \times n}$ has the form $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_p)$ where $p = \min(m, n)$. The diagonal values of Σ are called singular values of A and are usually ordered $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$.

We see that V, U contains the eigenvectors of $A^T A, A A^T$ as columns, respectively. Σ contains the square roots $\sqrt{\lambda_i}$ of the eigenvalues of $A^T A / A A^T$.

$$A^T A = V \Sigma^T \Sigma V^T = V \Sigma_n V^T$$

$$A A^T = U \Sigma^T \Sigma U^T = U \Sigma_m U^T$$

This way the SVD can also be used to diagonalize the symmetric positive semidefinite matrix $A^T A$ or $A A^T$.

It is often convenient to break the SVD into two parts. Split the singular values into non-zero and zero parts:

$$\Sigma = \begin{pmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{pmatrix} \quad \text{where} \quad \Sigma_1 = \text{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{r \times r} \quad \text{and} \quad \Sigma_2 = \text{diag}(0, \dots, 0) \in \mathbb{R}^{(m-r) \times (n-r)}$$

Then the SVD can be written (Moon, 1999):

$$A = (U_1 \quad U_2) \begin{pmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{pmatrix} \begin{pmatrix} V_1^T \\ V_2^T \end{pmatrix} = U_1 \Sigma_1 V_1^T = \sum_{i=1}^r \sigma_i \bar{u}_i \bar{v}_i^T$$

From this factorization, it is easy to determine the four fundamental subspaces of A :

- $\mathcal{R}(A) = \text{span}(U_1)$
- $\mathcal{N}(A) = \text{span}(U_2)$
- $\mathcal{R}(A^T) = \text{span}(V_1)$
- $\mathcal{N}(A^T) = \text{span}(V_2)$

Relating this back to our diagram, for the $m \times n$ matrix $A: \mathbb{R}^n \rightarrow \mathbb{R}^m$ the SVD factorizes into:

- U decomposes $\mathbb{R}^m = \mathcal{R}(A) \oplus \mathcal{N}(A^T) = \text{span}(U_1) \oplus \text{span}(U_2)$
- V decomposes $\mathbb{R}^n = \mathcal{R}(A^T) \oplus \mathcal{N}(A) = \text{span}(V_1) \oplus \text{span}(V_2)$

4.5 Pseudoinverse of linear operators

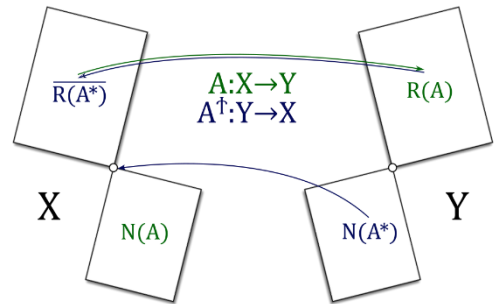
Definition: Best-approximate solution (Engl;Hanke;& Neubauer, 1996)

Let $T: \mathcal{X} \rightarrow \mathcal{Y}$ be a bounded linear operator where \mathcal{X} and \mathcal{Y} are Hilbert spaces .

- $x \in \mathcal{X}$ is called the least-squares solution of $Tx = y$ if $\|Tx - y\| = \inf\{\|Tz - y\| \mid z \in \mathcal{X}\}$
- $x \in \mathcal{X}$ is called best-approximate solution of $Tx = y$ if x is a least-squares solution of $Tx = y$ and $\|x\| = \inf\{\|z\| \mid z \text{ is a least squares solution of } Tx = y\}$

Bounded operator pseudoinverse (Engl;Hanke;& Neubauer, 1996)

Let $\bar{T} := T|_{\mathcal{N}(T)^\perp}: \mathcal{N}(T)^\perp \rightarrow \mathcal{R}(T)$. The Moore-Penrose inverse T^\dagger of $T \in \mathcal{L}(\mathcal{X}, \mathcal{Y})$ is defined as the unique linear extension of \bar{T}^{-1} to $\mathcal{D}(T^\dagger) = \mathcal{R}(T) \oplus \mathcal{R}(T)^\perp$ with $\mathcal{N}(T^\dagger) = \mathcal{R}(T)^\perp$.



Theorem (Engl;Hanke;& Neubauer, 1996)

Let $y \in \mathcal{D}(T^\dagger)$. Then $Tx = y$ has a unique best-approximative solution given by $x^\dagger = T^\dagger y$

where the set of all least-squares solutions is $x^\dagger + \mathcal{N}(T)$. $x \in \mathcal{X}$ is a least-squares solution of $Tx = y$ iff

$$T^*Tx = T^*y$$

The reason for assuming boundedness or $y \in \mathcal{D}(T^\dagger)$ is explained by two basic facts (Kreyszig, 1989):

- A linear operator $T: \mathcal{D}(T) \rightarrow Y$ is bounded iff it is continuous.
- If $\mathcal{R}(T)$ is closed, then $\mathcal{D}(T^\dagger) = \mathcal{R}(T) + \mathcal{R}(T)^\perp = \mathcal{Y}$.

And the following proposition (Engl;Hanke;& Neubauer, 1996):

- T^\dagger is bounded (continuous) if and only if $\mathcal{R}(T)$ is closed.

This implies that $\mathcal{R}(T)$ being closed is equivalent to the universal existence and continuous dependence of the solution. This is the case for example when T is finite dimensional. In general this is not the case. For example when a compact operator K is not finite dimensional and thus the range is not closed, the pseudoinverse K^\dagger is a densely defined unbounded linear operator. In this case the best-approximate solution of the inverse problem $Kx = y$ does not depend continuously on the right hand side and the problem is ill-posed. This implies means that further investigation is needed to stabilize the solution.

4.5.1 Compact operator pseudoinverse

For a compact operator, following condition implies the existence and spectral representation of the solution (Engl;Hanke;& Neubauer, 1996):

$$y \in \mathcal{D}(K^\dagger) \Leftrightarrow \sum_{i=1}^{\infty} \frac{|\langle y, u_i \rangle|^2}{\sigma_i^2} < \infty$$

$$\text{For } y \in \mathcal{D}(K^\dagger) : K^\dagger y = \sum_{i=1}^{\infty} \frac{\langle y, u_i \rangle}{\sigma_i} v_i$$

This expression also shows how the pseudoinverse is instable. Small errors in the data are amplified by large constants $\frac{1}{\sigma_i}$. Suppose for example we have added an error term $y_\delta = y + \delta u_i$. Then $\|y - y_\delta\| = \delta$ but $\|Ky - Ky_\delta\| = \left\| \frac{\langle \delta u_i, u_i \rangle}{\sigma_i} v_i \right\| = \frac{\delta}{\sigma_i} \rightarrow \infty$ as $i \rightarrow \infty$. Therefore if we assume that we have not observed y but a noisy sample y_δ close to it, $K^\dagger y_\delta \approx K^\dagger y$ is not a good solution. We seek to replace K^\dagger by a new operator K_α^\dagger which is continuous and stable against noise $K_\alpha^\dagger y_\delta \approx K^\dagger y$. This property is called regularization and the regularization parameter α determining the operator should be noise dependent so that as $\delta \rightarrow 0$ the solutions converge $K_\alpha^\dagger y_\delta \rightarrow K^\dagger y$.

4.5.2 Matrix pseudoinverse

Assume matrix $A \in \mathbb{R}^{m \times n}$ has the SVD:

$$A = U\Sigma V^T$$

where $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ are orthogonal and $\Sigma \in \mathbb{R}^{m \times n}$.

Then for the associated to the singular values $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_p)$ define the matrix:

$$\begin{aligned} \Sigma^\dagger &= \text{diag}(\sigma'_1, \sigma'_2, \dots, \sigma'_p) \\ \sigma'_i &= \begin{cases} 1/\sigma_i & \text{if } \sigma_i \neq 0 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

Based on the previous SVD factorization, it is easily seen that the pseudoinverse can be written:

$$A^\dagger = V\Sigma^\dagger U^T$$

4.6 Regularization of the pseudoinverse

Define the following notation for a compact operator $K^*Kx = \sum_{i=1}^{\infty} \sigma_i^2 \langle x, v_i \rangle v_i$:

$$f(K^*K) := \sum_{i=1}^{\infty} f(\sigma_i^2) \langle \cdot, v_i \rangle v_i$$

Then a regularized inverse is an operator which replaces the singular values $\sigma_i \hookrightarrow g_\alpha(\sigma_i)$ based on a function g_α (Engl; Hanke; & Neubauer, 1996):

$$x_\alpha^\dagger := K_\alpha^\dagger y = \sum_{i=1}^{\infty} g_\alpha(\sigma_i) \langle y, u_i \rangle v_i$$

where $g_\alpha: [0, \|T\|^2] \rightarrow \mathbb{R}$ is a piecewise continuous function such that

$$\begin{aligned} \forall \alpha > 0 \exists C > 0 : |\lambda g_\alpha(\lambda)| &\leq C \\ \lim_{\alpha \rightarrow 0} g_\alpha(\lambda) &= \frac{1}{\lambda} \end{aligned}$$

Then the following convergence result can be shown

$$\lim_{\alpha \rightarrow 0} g_\alpha(T^*T)T^*y = \begin{cases} T^\dagger y & \text{when } y \in \mathcal{D}(T^\dagger) \\ \infty & \text{when } y \notin \mathcal{D}(T^\dagger) \end{cases}$$

A recent application to machine learning is in (De Vito; Rosasco; & Verri, Spectral Methods for Regularization in Learning Theory, 2005) where different such functions are investigated in connection to algorithms.

Tikhonov regularization is a regularization method which defines the following g_α fulfilling the assumptions:

$$g_\alpha(\lambda) = \frac{1}{\lambda + \alpha}$$

Using the spectral theorem, then the regularized estimator can then be written through either T^*T or TT^* , and the solution expressed through the modified spectrum:

$$T_\alpha^\dagger y = \sum_{i=1}^{\infty} \frac{\sigma_n}{\sigma_n^2 + \alpha} \langle y, u_i \rangle v_i$$

$$T_\alpha^\dagger y = (T^*T + \alpha I)^{-1} T^* y = T^* (TT^* + \alpha I)^{-1} y$$

The Tikhonov regularized solution $x_\alpha^\dagger = T_\alpha^\dagger y$ is the unique minimizer of the Tikhonov functional (Engl; Hanke; & Neubauer, 1996) we saw previously:

$$f_\alpha(x) = \|Tx - y\|^2 + \alpha \|x\|^2$$

When the singular values are large compared to α , the factors $\frac{\sigma_n}{\sigma_n^2 + \alpha}$ are close to $\frac{1}{\sigma_n}$.

But as $n \rightarrow \infty$ the factors $\frac{\sigma_n}{\sigma_n^2 + \alpha} \rightarrow 0$ and they stay bounded; in particular when $\sigma_i < \alpha$ the factors are less than 1. This propagates the error $\langle y, u_i \rangle$ much less than the large factors $\frac{1}{\sigma_n}$, and therefore Tikhonov regularization stabilizes the solution as a function of the regularization strength α .

5 Theory of Kernel Approximation

In this section we present the mathematical formulation of learning theory in connection to kernel approximation. Operator theoretic results relevant to the approximation and our learning bound are also briefly discussed.

5.1 Random variable formulation

Define the learning scenario as follows. Assume the input space X is a closed and bounded subset of \mathbb{R}^n . Let the output space Y be a closed interval $[-M, M] \subseteq \mathbb{R}$. The sample space is defined $Z = X \times Y$. This is a common practical setting in machine learning (De Vito; Rosasco; & Verri, Spectral Methods for Regularization in Learning Theory, 2005), but it could be further generalized by assuming X to be a manifold and $Y = \mathbb{R}^k$ (Cucker & Smale, 2001).

With the sample space we have associated an unknown probability distribution expressing the target:

$$\rho(x, y) = \rho(y|x)\rho_X(x)$$

This formulation encompasses the scenario where target function f^* has a random noise term σ :

$$\rho(x, y) = f^*(x) + \sigma(x, y)$$

The training set \bar{z} is then drawn i.i.d from Z according to ρ :

$$\bar{z} = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

To study the learning setting theoretically, define the **empirical error** $\mathcal{E}_{\bar{z}}: \mathcal{H} \rightarrow \mathbb{R}$ and **population error** $\mathcal{E}: \mathcal{H} \rightarrow \mathbb{R}$:

$$\mathcal{E}_{\bar{z}}(f) = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)^2$$

$$\mathcal{E}(f) = \int_{X \times Y} (f(x) - y)^2 d\rho(x, y)$$

We are interested in a function f which minimizes the error $\mathcal{E}(f)$ over the space of square integrable functions $L^2(X, p_X)$. For this purpose, define the regression function as

$$f_\rho(x) = \int_Y y d\rho(y|x)$$

Lemma (Cucker & Smale, 2001)

The error for $f \in L^2(X, p_X)$ can be decomposed:

$$\mathcal{E}(f) = \int_X \left(f(x) - f_\rho(x) \right)^2 d\rho_X(x) + \mathcal{E}(f_\rho)$$

Proof

$$\begin{aligned} \mathcal{E}(f) &= \int_Z (f(x) - f_\rho(x) + f_\rho(x) - y)^2 d\rho(x, y) \\ &= \int_X \left(f(x) - f_\rho(x) \right)^2 d\rho_X(x) + \int_Z (f_\rho(x) - y)^2 d\rho(x, y) \\ &\quad + 2 \int_Z (f(x) - f_\rho(x))(f_\rho(x) - y) d\rho(x, y) \\ &= \int_X \left(f(x) - f_\rho(x) \right)^2 d\rho_X(x) + \mathcal{E}(f_\rho) \end{aligned}$$

The error $\mathcal{E}(f)$ of function f decomposes into two terms. The second term is the error associated with the regression function f_ρ , is independent of f , and the first term is the expected squared difference between f and f_ρ . This also shows that the regression function achieves the minimum error and that in considering the minimizer of $\mathcal{E}(f)$, the term $\mathcal{E}(f_\rho)$ may be omitted since it is independent of f .

However, since the probability distribution $\rho(x, y)$ is unknown, we cannot use it to define the regression function, which is thus also unknown. We would like to approximate the regression function as closely as possible, given the sample \bar{z} that we do know and the hypothesis space \mathcal{H} we have. The goal is thus to find a learning algorithm \mathcal{L} :

$$\mathcal{L}: \bar{z} \mapsto f_{\bar{z}} \text{ such that } f_{\bar{z}} \text{ approximates } f_\rho$$

where the subscript highlights the fact that the function is deduced from the sample \bar{z} .

The learning algorithm implicitly or explicitly is based on a family of functions \mathcal{H} used for learning. In general it may be that $f_\rho \notin \mathcal{H}$. Given a fixed hypothesis set, in that case our goal is to learn $f_{\mathcal{H}}$ defined as the closest approximation of f_ρ in \mathcal{H} . Since

this is unavailable, we have also define a hypothesis $f_{\bar{z}}$ which minimizes the empirical error:

$$f_{\mathcal{H}} = \operatorname{argmin}_{f \in \mathcal{H}} \mathcal{E}(f)$$

$$f_{\bar{z}} = \operatorname{argmin}_{f \in \mathcal{H}} \mathcal{E}_{\bar{z}}(f)$$

Define the error of a function $f \in \mathcal{H}$ as:

$$\mathcal{E}_{\mathcal{H}}(f) = \mathcal{E}(f) - \mathcal{E}(f_{\mathcal{H}})$$

This allows the following decomposition for a function $f_{\bar{z}}$ (Cucker & Smale, 2001)

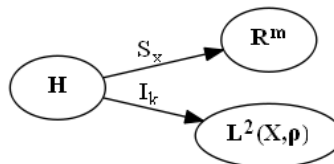
$$\mathcal{E}(f_{\bar{z}}) = \mathcal{E}_{\mathcal{H}}(f_{\bar{z}}) + \mathcal{E}(f_{\mathcal{H}})$$

The first term $\mathcal{E}_{\mathcal{H}}(f_{\bar{z}})$ is the error incurred due to using $f_{\bar{z}}$ instead of the best available function in \mathcal{H} , called the *sampling error*, and the second term is the error incurred due to using \mathcal{H} to approximate f_{ρ} , called the *approximation error*. Theoretical considerations often analyse these errors separately.

Algorithms that aim to minimize a population error $\mathcal{E}(f)$ based on the empirical error $\mathcal{E}_{\bar{z}}(f)$, with a possibly added regularization functional also known as complexity penalty, are called **empirical risk minimization** algorithms. The least squares algorithm is a special case of such an algorithm in the sense that the loss is not only convex but can be defined through a norm, which defines a convex loss by the triangle inequality. We now analyse the minimizers and their properties by contrasting them based on this fact.

5.2 Abstract minimizers

Define $\|x\|_m = \frac{1}{m} \|x\|_{\mathbb{R}^m}$, a norm on the sample output space $y_i \in \mathbb{R}^m$, and $\|f\|_{\rho} = \left(\int_{X \times Y} f(x)^2 d\rho(x) \right)^{\frac{1}{2}}$ where $\rho(x)$ is a shorthand for the marginal measure $\rho_X(x) = \int_Y \rho(x, y)$, which defines a norm on the space of square integrable functions $L^2(X, \rho)$. Define the sampling operator $S_{\bar{x}}: \mathcal{H} \rightarrow \mathbb{R}^m$ as before and the inclusion $I_k: \mathcal{H} \rightarrow L^2(X, \rho)$ as the map of f to its equivalence class in $L^2(X, \rho)$:



$$S_{\bar{x}}f = (f(x_1), \dots, f(x_m))^T$$

$$I_k f = [f]$$

The inclusion operator I_k is associated with a number of technical considerations:

- If ρ is degenerate, $I_k f$ is an equivalence class of functions and evaluation $I_k f(x)$ is not well-defined (Steinwart & Christmann, 2008). For simplicity some theoretical literature assumes that ρ is non-degenerate, to avoid considerations involving zero-measure sets. Then the inclusion is injective, because $I_k f = 0$ implies that $f(x) \neq 0$ is possible in some open set of ρ -measure zero, which must be empty (ρ nondegenerate). This implies if k is continuous $f(x) = 0$, otherwise identically zero for ρ almost all $x \in X$. Thus $(I_k f)(x) = \langle f, K_x \rangle_H = f(x)$ is well-defined. In this case I_k merely changes the norm from $\|f\|_H$ to $\|f\|_{L^2}$.
- The range of inclusion operator needs to be a subset of $L^2(X, \rho)$. When the kernel is measurable and bounded $I_k(\mathcal{H}) \subset L^2(X, \rho)$. Since $k(x, y)^2 \leq k(x, x)k(y, y)$, sufficient condition for boundedness is that the kernel is continuous and X is compact, because then the kernel is bounded:
$$\sup_{x \in X} \sqrt{k(x, x)} = \sup_{x \in X} \|K_x\|_H < \infty$$
- The range $I_k(\mathcal{H})$ is in some theoretical considerations required to be dense in $L^2(X, \rho)$ so that $\inf_{f \in \mathcal{H}} \mathcal{E}(I_k f) = \inf_{f \in L^2(X, \rho)} \mathcal{E}(f) = \mathcal{E}(f_\rho)$. However, if this is not the case, the optimal solution is the projection of f_ρ into the closure of $I_k(\mathcal{H})$, often denoted by Pf_ρ (De Vito; Rosasco; & Verri, Spectral Methods for Regularization in Learning Theory, 2005).

Based on our discussion of subspaces associated to a linear operator and basic properties (Kreyszig, 1989) of the adjoint, I_k^* , which to be seen has an informative form for kernels, we have on injectivity/surjectivity:

- $\mathcal{R}(I_k)$ is dense in $L^2(X, \rho) \Leftrightarrow \mathcal{N}(I_k^*) = \{0\}$
- $\mathcal{N}(I_k) = \{0\} \Leftrightarrow \mathcal{R}(I_k^*)$ is dense in \mathcal{H}

Then the errors may be written through the norms:

$$\mathcal{E}_{\bar{z}}(f) = \|S_{\bar{x}}f - \bar{y}\|_m^2$$

$$\mathcal{E}(f) = \|I_k f - y\|_\rho$$

The augmented empirical and population error add the norm penalty $\lambda > 0$:

$$\mathcal{E}_{\bar{z}}^{\lambda}(f) = \|S_{\bar{x}}f - \bar{y}\|_m^2 + \lambda\|f\|_H^2$$

$$\mathcal{E}^{\lambda}(f) = \|I_k f - y\|_{\rho} + \lambda\|f\|_H^2$$

Notice that even though we seek solution in \mathcal{H} , the error is defined through the norm in $L^2(X, \rho)$, so we need to combine the function with inclusion I_k for this to be defined. Using the norm formulation and the two linear operators $S_{\bar{x}}$ and I_k between Hilbert spaces, the minimization problems become linear inverse problems. As was seen in the general inverse case, for an expressive hypothesis space regularization is necessary to enforce stability since generalization is of central importance in learning. The solution to the regularized empirical and population least squares error is defined:

$$f_{\bar{z}}^{\lambda} = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \|S_{\bar{x}}f - \bar{y}\|_n^2 + \lambda\|f\|_H^2$$

$$f^{\lambda} = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \|I_k f - f_{\rho}\|_{\rho} + \lambda\|f\|_H^2$$

The unregularized solutions to the empirical and population cases can be defined as the pseudoinverse solutions to the linear inverse problems (Cucker & Smale, 2001):

$$S_{\bar{x}}f = \bar{y}$$

$$I_k f = f_p$$

The regularized solutions can be obtained for example through the functional derivative (Smale & Zhou, 2007) or modifying the singular values in the pseudoinverse series expansion if the operator is compact. Because adding λI ($\lambda > 0$) to a self-adjoint operator makes the eigenvalues strictly positive, these have unique solutions expressible through the approximator operator (De Vito; Rosasco; Caponetto; De Giovannini; & Odone, 2005):

$$f_{\bar{z}}^{\lambda} = (S_{\bar{x}}^* S_{\bar{x}} + \lambda I)^{-1} S_{\bar{x}}^* \bar{y}$$

$$f^{\lambda} = (I_K^* I_K + \lambda I)^{-1} I_K^* f_p$$

It was previously stated that as $\lambda \rightarrow 0$ these operators converge to the pseudo-inverse.

5.3 Kernel minimizers

Assume that \mathcal{H} is separable and for the kernel k :

$$\|k\|_\rho = \left(\int_X k(x, x) d\rho(x) \right)^{\frac{1}{2}} < \infty$$

This requirement is satisfied if $k: X \times X \rightarrow \mathbb{R}$ is continuous on a compact set X .

It is immediately verified that $I_K: \mathcal{H} \rightarrow L^2(X, \rho)$ is continuous (bounded) (Steinwart & Christmann, 2008):

$$\int_X |f(x)|^2 d\rho(x) = \int_X |\langle f, k(\cdot, x) \rangle|^2 d\rho(x) \leq \|f\|_H^2 \int_X k(x, x) d\rho(x)$$

This implies that the adjoint I_K^* exists and is likewise bounded since $\|I_K^*\| = \|I_K\|$ (Kreyszig, 1989).

The adjoint $S_{\bar{x}}^*: \mathbb{R}^n \rightarrow \mathcal{H}$ is given by (Rosasco; Belkin; & De Vito, 2010):

$$S_{\bar{x}}^* \bar{c} = \frac{1}{m} \sum_{i=1}^m c_i K_{x_i}$$

This follows trivially from the definition $\langle S_{\bar{x}} f, \bar{c} \rangle_{\mathbb{R}^m} = \frac{1}{m} \sum_{i=1}^m \langle f, K_{x_i} \rangle_H c_i = \left\langle f, \frac{1}{m} \sum_{i=1}^m c_i K_{x_i} \right\rangle_H$.

The adjoint $I_K^*: L^2(X, \rho) \rightarrow \mathcal{H}$ is given by (Rosasco; Belkin; & De Vito, 2010):

$$(I_K^* f)(x) = \int_X k(x, y) f(y) d\rho(y)$$

This follows from the definition $\langle I_K^* f, g \rangle_H = \langle f, I_K g \rangle_{L^2}$. Setting $g = k(\cdot, x)$, we have $(I_K^* f)(x) = \langle f, I_K k(\cdot, x) \rangle_{L^2} = \int_X k(x, y) f(y) d\rho(y)$.

Therefore we may define the following **empirical operators**, the **kernel matrix** $K: \mathbb{R}^m \rightarrow \mathbb{R}^m$ and the **covariance operator** $C: \mathcal{H} \rightarrow \mathcal{H}$:

$$K = S_{\bar{x}} S_{\bar{x}}^*$$

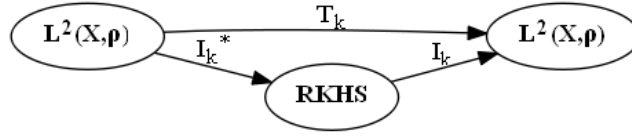
$$C = S_{\bar{x}}^* S_{\bar{x}}$$

The following expressions are derived for reproducing kernel Hilbert spaces:

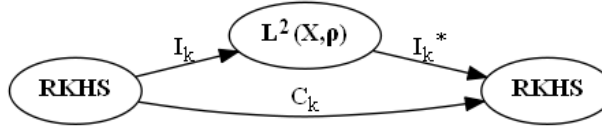
$$K\bar{c} = S_{\bar{x}}S_{\bar{x}}^*\bar{c} = \begin{pmatrix} \langle \frac{1}{m} \sum_{i=1}^m c_i K_{x_i}, K_{x_1} \rangle \\ \dots \\ \langle \frac{1}{m} \sum_{i=1}^m c_i K_{x_i}, K_{x_n} \rangle \end{pmatrix} \bar{c} = \frac{1}{m} \begin{pmatrix} k(x_1, x_1) & \dots & k(x_n, x_1) \\ \vdots & \ddots & \vdots \\ k(x_1, x_n) & \dots & k(x_n, x_n) \end{pmatrix} \bar{c}$$

$$Cf = S_{\bar{x}}^*S_{\bar{x}}f = S_{\bar{x}}^* \begin{pmatrix} f(x_1) \\ \dots \\ f(x_n) \end{pmatrix} = \frac{1}{m} \sum_{i=1}^m f(x_i) K_{x_i}$$

Also define following **population operators**, the **kernel operator** $T_k: L^2(X, \rho) \rightarrow L^2(X, \rho)$ and the **covariance operator** $C_k: \mathcal{H} \rightarrow \mathcal{H}$ (Hein & Olivier, 2004):



$$T_k = I_k I_k^*$$



$$C_k = I_k^* I_k$$

Using $(I_k f)(x) = \langle f, k(\cdot, x) \rangle_H$ and the fact that the integral commutes with the inner product:

$$T_k f = I_k I_k^* f = \left\langle \int_X K_y f(y) d\rho_X(y), K_x \right\rangle_H = \int_X k(x, y) f(y) d\rho(y)$$

$$C_k f = I_k^* I_k f = \int_X K_y (I_k f)(y) d\rho_X(y) = \int_X K_y \langle f, K_y \rangle_H d\rho(y)$$

The first integral converges in $L^2(X, \rho)$ norm and the second in \mathcal{H} norm.

The population operators closely connected to their empirical variants, for example the following descriptions are used in the literature:

	Empirical	Population
Kernel	$\frac{1}{m} S_{\bar{x}} S_{\bar{x}}^*$	$I_k I_k^*$

Covariance	$\frac{1}{m} S_{\bar{x}}^* S_{\bar{x}}$	$I_k^* I_k$
-------------------	---	-------------

Therefore the unique solution $f_{\bar{z}}^\lambda = \left(\frac{1}{m} S_{\bar{x}}^* S_{\bar{x}} + \lambda I \right)^{-1} \frac{1}{m} S_{\bar{x}}^* \bar{y} = \frac{1}{m} S_{\bar{x}}^* \left(\frac{1}{m} S_{\bar{x}} S_{\bar{x}}^* + \lambda I \right)^{-1} \bar{y}$ in \mathcal{H} minimizing the regularized least squares error can be written through either the kernel matrix or the covariance operator. Using the kernel matrix, we obtain the following finite dimensional problem:

$$f_{\bar{z}}^\lambda(x) = \sum_{i=1}^m a_i K(x, x_i) \text{ where } \bar{a} = \left(\frac{1}{m} K + \lambda I \right)^{-1} \frac{1}{m} \bar{y}$$

Using the previous result that $(I_k^* I_k + \lambda I)^{-1} I_k^* = I_k^* (I_k I_k^* + \lambda I)^{-1}$, the population error minimizer can be written through either T_k or C_k and optionally combined with the inclusion I_k , the resulting function considered in either \mathcal{H} or $L^2(X, \rho)$. The population error minimizer as an element in $L^2(X, \rho)$ is given by:

$$f^\lambda = (T_k + \lambda I)^{-1} T_k f_\rho$$

The specific form used in the literature is quite variable, contrast for example (Smale & Zhou, 2007) (Rosasco;Belkin;& De Vito, 2010) (De Vito;Rosasco;Caponetto;De Giovannini;& Odone, 2005).

5.4 Spectral and singular value decompositions

The following stronger results, contingent on $\|A\| \leq \|A\|_{\text{HS}} \leq \|A\|_{\text{TC}}$ (Rosasco;Belkin;& De Vito, 2010), may be established under the assumptions. See the appendix (Steinwart & Christmann, 2008, ss. 507-529) which has additional straightforward results founded on these operator norms, closely connected to the spectral theorem in (Kreyszig, 1989). These results allow intuitive interpretations for the operators.

Lemma: The inclusion I_k and adjoint I_k^* are Hilbert-Schmidt and thus compact (Steinwart & Christmann, 2008). For ONB $(e_i)_{i \geq 1}$ of \mathcal{H} :

$$\|I_k^*\|_{\text{HS}} = \sum_{i=1}^{\infty} \|I_k^* e_i\|_\rho^2 = \int_X \sum_{i=1}^{\infty} |I_k^* e_i(x)|^2 d\rho(x) = \int_X \sum_{i=1}^{\infty} e_i(x)^2 d\rho(x) = \|k\|_\rho < \infty$$

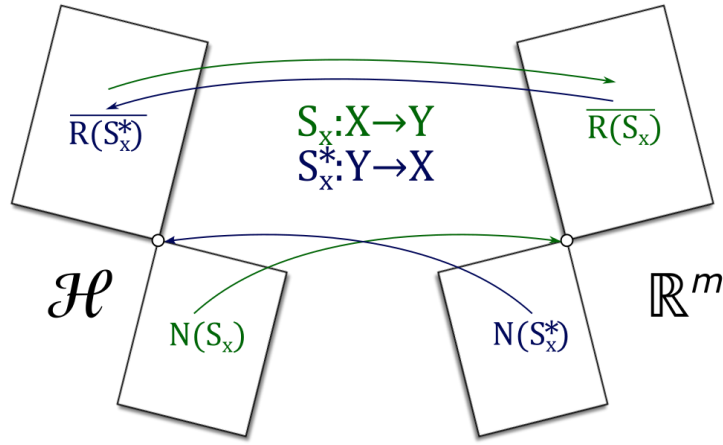
For I_k^* , which implies the same for I_k since $\|A\|_{\text{HS}} = \|A^*\|_{\text{HS}}$.

This leads us to the following theorem, where compactness is essential to be able to apply several results in functional analysis.

Theorem: $T_k = I_k I_k^*$ and $C_k = I_k^* I_k$ are compact, positive, self-adjoint and nuclear operators with $\|T_k\|_{\text{nuc}} = \|C_k\|_{\text{nuc}} = \|I_k^*\|_{\text{HS}}$.

Proof: Compactness of I_k^* implies compactness of $I_k I_k^*$ and $\|I_k^*\|_{\text{HS}} < \infty \Rightarrow \|I_k I_k^*\|_{\text{nuc}} = \|I_k^*\|_{\text{HS}}$ (Steinwart & Christmann, 2008). Case for C_k goes analogously.

5.4.1 Empirical decomposition



Factor $\mathcal{H} = \mathcal{N}(S_{\bar{x}}) \otimes \overline{\mathcal{R}(S_{\bar{x}}^*)}$ and $\mathbb{R}^m = \mathcal{N}(S_{\bar{x}}^*) \oplus \overline{\mathcal{R}(S_{\bar{x}})}$.

Then we have the following factorizations for operators $S_{\bar{x}}: \mathcal{H} \rightarrow \mathbb{R}^m, S_{\bar{x}}^*: \mathbb{R}^m \rightarrow \mathcal{H}, S_{\bar{x}} S_{\bar{x}}^*: \mathbb{R}^m \rightarrow \mathbb{R}^m, S_{\bar{x}}^* S_{\bar{x}}: \mathcal{H} \rightarrow \mathcal{H}$:

$S_{\bar{x}} f = \sqrt{\lambda_i} \langle f, \hat{v}_i \rangle x_i$	$Kx := S_{\bar{x}} S_{\bar{x}}^* x = \lambda_i \langle x, x_i \rangle x_i$
$S_{\bar{x}}^* x = \sqrt{\lambda_i} \langle x, x_i \rangle \hat{v}_i$	$Cf := S_{\bar{x}}^* S_{\bar{x}} f = \lambda_i \langle f, \hat{v}_i \rangle \hat{v}_i$

Where $(\{x_i\}_{i \geq 1}, \{\hat{v}_i\}_{i \geq 1}, \{\lambda_i\}_{i \geq 1})$ is a singular system of orthonormal eigenvectors of K , orthonormal eigenvectors of C and nonzero eigenvalues of K/C . The sets $\{x_i\}_{i \geq 1}, \{\hat{v}_i\}_{i \geq 1}$ are also a basis for $\mathcal{R}(S_{\bar{x}}), \mathcal{R}(S_{\bar{x}}^*)$, respectively.

The following kernel specific explicit expressions can be written:

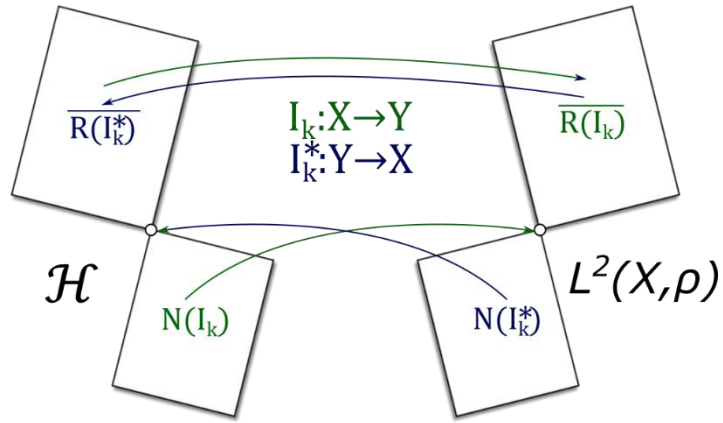
$$S_{\bar{x}} \hat{v}_i = \sqrt{\lambda_i} x_i \Rightarrow x_i = \frac{1}{\sqrt{\lambda_i}} (\hat{v}_i(x_1), \dots, \hat{v}_i(x_m))^T \in \mathbb{R}^m$$

$$S_x^* x_i = \sqrt{\lambda_i} \hat{v}_i \Rightarrow \hat{v}_i = \frac{1}{\sqrt{\lambda_i} m} \sum_{i=1}^m x_i K_{x_i} \in \mathcal{H}$$

The latter expression is sometimes called **kernel PCA** or **KPCA**, because it can be seen as an infinite dimensional analog to the PCA of a finite dimensional feature space. Given a data point x and i 'th eigenvector x_i of K , the projection to i 'th component may be computed for each \hat{v}_i :

$$\hat{v}_i(x) = \frac{1}{\sqrt{\lambda_i} m} \sum_{i=1}^m x_i k(x, x_i)$$

5.4.2 Population decomposition



Factor $\mathcal{H} = \mathcal{N}(I_k) \otimes \overline{\mathcal{R}(I_k^*)}$ and $L^2(X, \rho) = \mathcal{N}(I_k^*) \oplus \overline{\mathcal{R}(I_k)}$.

Then we have the following factorizations for operators $I_k: \mathcal{H} \rightarrow L^2(X, \rho)$, $I_k^*: L^2(X, \rho) \rightarrow \mathcal{H}$, $I_k I_k^*: L^2(X, \rho) \rightarrow L^2(X, \rho)$, $I_k^* I_k: \mathcal{H} \rightarrow \mathcal{H}$:

$I_k f = \sqrt{\mu_i} \langle f, v_i \rangle u_i$	$T_k g := I_k I_k^* g = \mu_i \langle g, u_i \rangle u_i$
$I_k^* g = \sqrt{\mu_i} \langle g, u_i \rangle v_i$	$C_k f := I_k^* I_k f = \mu_i \langle f, v_i \rangle v_i$

Where $(\{u_i\}_{i \geq 1}, \{v_i\}_{i \geq 1}, \{\mu_i\}_{i \geq 1})$ is a singular system of orthonormal eigenvectors of T_k , orthonormal eigenvectors of C_k and nonzero eigenvalues of T_k/C_k . The sets $\{u_i\}_{i \geq 1}, \{v_i\}_{i \geq 1}$ are also a basis for $\mathcal{R}(I_k), \mathcal{R}(I_k^*)$, respectively.

The following kernel specific explicit expressions can be written:

$$I_k v_i = \sqrt{\mu_i} u_i \Rightarrow u_i(y) = \frac{1}{\sqrt{\mu_i}} v_i(x) \text{ for } \rho\text{-almost all } x \in X$$

$$I_k^* u_i = \sqrt{\mu_i} v_i \Rightarrow v_i(x) = \frac{1}{\sqrt{\mu_i}} \int_{y \in X} k(x, y) u(y) \rho(y) \in \mathcal{H}$$

5.4.3 Decompositions and associated projections

Consider the following operators:

$$\bar{y}_{emp} = S_{\bar{x}} S_{\bar{x}, \lambda}^\dagger \bar{y} = S_{\bar{x}} S_{\bar{x}}^* (S_{\bar{x}} S_{\bar{x}}^* + \lambda I)^{-1} \bar{y} = K(K + \lambda I)^{-1} \bar{y}$$

$$f_{emp}^H = S_{\bar{x}, \lambda}^\dagger S_{\bar{x}} f = (S_{\bar{x}}^* S_{\bar{x}} + \lambda I)^{-1} S_{\bar{x}}^* S_{\bar{x}} f = (C + \lambda I)^{-1} C f$$

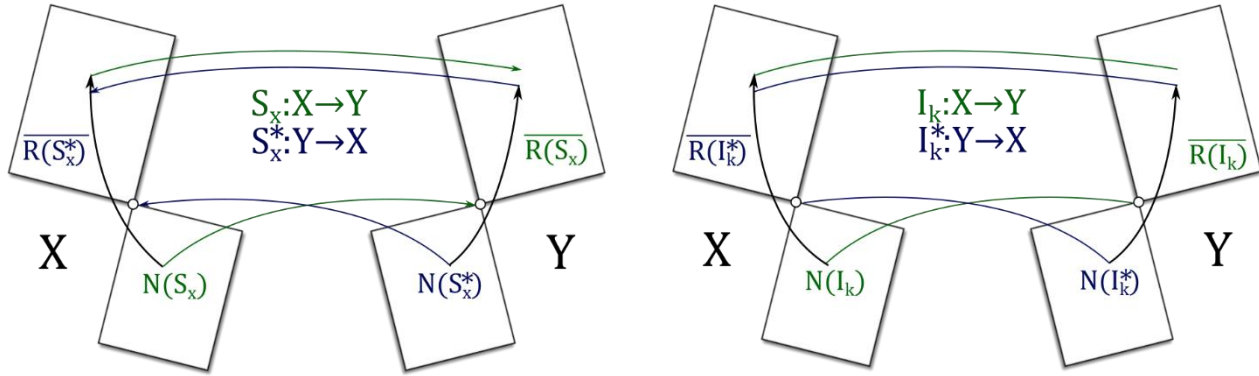
$$g_{pop} = I_k I_{k, \lambda}^\dagger g = I_k I_k^* (I_k I_k^* + \lambda I)^{-1} g = T_k (T_k + \lambda I)^{-1} g$$

$$f_{pop}^H = I_{k, \lambda}^\dagger I_k f = (I_k^* I_k + \lambda I)^{-1} I_k^* I_k f = (C_k + \lambda I)^{-1} C_k f$$

These can be interpreted as follows:

- Given data \bar{y} , predictions \bar{y}_{emp} given by a function learned using \mathcal{H} .
- Given function $f \in \mathcal{H}$, function f_{emp}^H learned from $f(\bar{x})$ using \mathcal{H} .
- Given function $g \in L^2(X, \rho)$, function in $g_{pop} \in L^2(X, \rho)$ learned using \mathcal{H} .
- Given function $f \in \mathcal{H}$, function f_{pop}^H learned from $I_k f \in L^2(X, \rho)$ using \mathcal{H} .

Consider the spectral representations of the operators and the following projections:



$K(K + \lambda I)^{-1} \bar{y}$ $= \sum_{i \geq 1} \frac{\lambda_i}{\lambda_i + \lambda} \langle x, x_i \rangle_{\mathbb{R}^m} x_i$	$P_{\overline{\mathcal{R}(S_{\bar{x}})}} \bar{y} = \sum_{i \geq 1: \lambda_i > 0} \langle x, x_i \rangle_{\mathbb{R}^m} x_i$
$(C + \lambda I)^{-1} C f = \sum_{i \geq 1} \frac{\lambda_i}{\lambda_i + \lambda} \langle f, \hat{v}_i \rangle_H \hat{v}_i$	$P_{\overline{\mathcal{R}(S_{\bar{x}}^*)}} f = \sum_{i \geq 1: \lambda_i > 0} \langle f, \hat{v}_i \rangle_H \hat{v}_i$

$T_k(T_k + \lambda I)^{-1}g = \sum_{i \geq 1} \frac{\mu_i}{\mu_i + \lambda} \langle g, u_i \rangle u_i$	$P_{\overline{\mathcal{R}(I_k)}}g = \sum_{i \geq 1: \mu_i > 0} \langle g, u_i \rangle u_i$
$(C_k + \lambda I)^{-1}C_k f = \sum_{i \geq 1} \frac{\mu_i}{\mu_i + \lambda} \langle f, v_i \rangle v_i$	$P_{\overline{\mathcal{R}(I_k^*)}}f = \sum_{i \geq 1: \mu_i > 0} \langle f, v_i \rangle_H v_i$

As $\lambda \rightarrow 0$, the regularized pseudoinverse converges to the pseudoinverse, and these operators converge to the corresponding projections. They can be interpreted as ‘filtered projections’, in the sense that the term λ scales back the projection coefficients based on the regularized pseudoinverse filtering.

5.5 Connection between \mathcal{H} and L^2 and Mercer’s theorem

5.5.1 RKHS isometrically embedded in L^2

Consider a measure μ on $L^2(X, \mu)$. Define a function $f: \text{supp}(\mu) \rightarrow \mathbb{R}$ associated to each equivalence class $[(f)]$ in $L^2(X, \mu)$, which is then well defined:

$$f := g|_{\text{supp}(\mu)} \text{ for any representer } g \in [(f)]$$

Let $\text{supp}(\mu) = X$, i.e. $\mu(X') > 0$ for every open subset $X' \subset X$. It was previously stated for a continuous kernel $k: X \times X \rightarrow \mathbb{R}$ the inclusion $I_k: \mathcal{H} \rightarrow L^2(X, \mu)$ is then injective, and if X is compact the functions in RKHS are bounded so inclusion is well defined. Consider the decomposition from before:

$$I_k e_i = e_i = \sqrt{\lambda_i} \tilde{e}_i \in L^2(X, \mu)$$

Where $\{\tilde{e}_i\}_{i \geq 1} \subset L^2(X, \mu)$ is an ONB of $\mathcal{N}(I_k^*)^\perp$ and $\{e_i\}_{i \geq 1} \subset \mathcal{H}$ is an ONB of $\mathcal{N}(I_k)^\perp = \mathcal{H}$, $\{\lambda_i\}_{i \geq 0}$ are the eigenvalues of $C_k = I_k I_k^*$ / the nonzero eigenvalues of $T_k = I_k^* I_k$.

Inclusion maps to the same function by definition $e_i(x) = \sqrt{\lambda_i} \tilde{e}_i(x) \forall x \in X$, but changes the norms: if $f = \sum_{i \geq 1} \alpha_i e_i$ then $I_k(f) = \sum_{i \geq 1} \sqrt{\lambda_i} \alpha_i \tilde{e}_i$, implying $\|f\|_H = \sum_{i \geq 1} \alpha_i^2$ and $\|I_k(f)\|_{L^2} = \sum_{i \geq 1} \lambda_i \alpha_i^2$. Since inclusion is injective, I_k is a bijection between \mathcal{H} and $\mathcal{N}(I_k^*)^\perp$. We may therefore consider the ‘reverse inclusion’:

$$I_k^{-1} \tilde{e}_i = \tilde{e}_i = \frac{1}{\sqrt{\lambda_i}} e_i \in \mathcal{H}$$

Now consider the following map \mathfrak{S} between the bases of $\mathcal{N}(I_k^*)^\perp$ and \mathcal{H} :

$$\mathfrak{I}\tilde{e}_i := e_i = \sqrt{\lambda_i}\tilde{e}_i$$

This maps to a different function, but retains norms: if $f = \sum_{i \geq 1} \beta_i \tilde{e}_i$ then $\mathfrak{I}(f) = \sum_{i \geq 1} \sqrt{\lambda_i} \beta_i \tilde{e}_i = \sum_{i \geq 1} \beta_i e_i$, implying $\|f\|_{L^2} = \sum_{i \geq 1} \beta_i^2$ and $\|\mathfrak{I}(f)\|_H = \sum_{i \geq 1} \beta_i^2$. It is clearly a bijection between $\mathcal{N}(I_k^*)^\perp$ and \mathcal{H} , inducing an isometric isomorphism between the spaces.

Since T_k is a positive self-adjoint compact operator, the square root $T_k = T_k^{\frac{1}{2}} \circ T_k^{\frac{1}{2}}$ is defined and given by:

$$T_k^{\frac{1}{2}} f = \sum_{i \geq 1} \sqrt{\lambda_i} \langle f, \tilde{e}_i \rangle_{L^2} \tilde{e}_i$$

Considering $T_k^{\frac{1}{2}} f$ as an element of \mathcal{H} , we see that it coincides with \mathfrak{I} and induces an isometric isomorphism between $\mathcal{R}(I_k) = \mathcal{N}(I_k^*)^\perp$ and \mathcal{H} .

5.5.2 Mercer representation

Recall that $k: X \times X \rightarrow \mathbb{R}$ was assumed continuous, X compact and $\text{supp}(\mu) = X$.

We know from before that if $\{e_i\}_{i \geq 1}$ is an ONB of \mathcal{H} , then

$$k(x, y) = \sum_{i \geq 1} e_i(x) e_i(y)$$

Where the convergence is uniform and absolute. However, since $e_i(x) = \sqrt{\lambda_i} \tilde{e}_i(x)$, we can express $k(x, y)$ in terms of the eigenfunctions $\{\tilde{e}_i\}_{i \geq 1}$ associated to nonzero eigenvalues of the operator $T_k = I_k^* I_k$:

$$k(x, y) = \sum_{i \geq 1} \lambda_i \tilde{e}_i(x) \tilde{e}_i(y)$$

This result is known as **Mercer's theorem** in the literature. It can be used to derive an expression for k , and an ONB representation of the RKHS, having obtained the eigenfunctions of the positive self-adjoint compact operator $T_k: L^2(X, \mu) \rightarrow L^2(X, \mu)$.

Since $k(x, y) = \langle \sqrt{\lambda_i} \tilde{e}_i(x), \sqrt{\lambda_i} \tilde{e}_i(y) \rangle_{l^2}$, Mercer's theorem also gives a feature map $\phi: X \rightarrow l^2(J)$ related to a spectral decomposition (Steinwart & Christmann, 2008):

$$\phi: x \mapsto \left(\sqrt{\lambda_i} \tilde{e}_i(x) \right)_{i \geq 1}$$

Which is well defined because $\sum_{i \geq 1} |\sqrt{\lambda_i} e_i(x)|^2 = k(x, x) < \infty$

5.5.3 ONB construction of a RKHS

Under the assumptions of Mercer's theorem, an alternative construction exists for the RKHS in terms of eigenfunctions $\{\tilde{e}_i\}_{i \geq 1} \subset L^2(X, \mu)$ associated to nonzero eigenvalues $\{\lambda_i\}_{i \geq 1}$ of T_k . Define (Steinwart & Christmann, 2008):

$$\mathcal{H} = \left\{ f = \sum_{j \in J} \sqrt{\lambda_j} a_j \tilde{e}_j : \{a_j\} \in l^2(J) \right\}$$

With an inner product

$$\left\langle \sum_{j \in J} \sqrt{\lambda_j} a_j \tilde{e}_j, \sum_{j \in J} \sqrt{\lambda_j} b_j \tilde{e}_j \right\rangle = \sum_{j \in J} a_j b_j$$

Then \mathcal{H} is the RKHS of k and $\{\sqrt{\lambda_i} \tilde{e}_i\}_{i \geq 1}$ clearly an ONB of \mathcal{H} .

Proof

The functions are well-defined since they converge for all x , from Hölder's inequality (Kreyszig, 1989, s. 14):

$$\sum_{j \in J} |\sqrt{\lambda_j} a_j \tilde{e}_j| \leq \left[\sum_{j \in J} |a_j|^2 \right]^{\frac{1}{2}} \left[\sum_{j \in J} |\sqrt{\lambda_j} \tilde{e}_j(x)|^2 \right]^{\frac{1}{2}} = \|\{a_j\}\|_{l^2} \sqrt{k(x, x)}$$

Then k obtained from Mercer's theorem is the reproducing kernel of \mathcal{H} :

- $k(\cdot, x) \in \mathcal{H}$, since $k(\cdot, x) = \sum_{j \in J} (\lambda_j \tilde{e}_j(x)) \tilde{e}_j \Rightarrow \sum_{j \in J} |\sqrt{\lambda_j} \tilde{e}_j(x)|^2 = \sum_{j \in J} \lambda_j \tilde{e}_j(x)^2 = k(x, x) < \infty$
- k has the reproducing property since $f = \sum_{j \in J} a_j \tilde{e}_j \Rightarrow \langle f, k(\cdot, x) \rangle = \langle \sum_{j \in J} a_j \tilde{e}_j, \sum_{j \in J} (\lambda_j \tilde{e}_j(x)) \tilde{e}_j \rangle = \sum_{j \in J} a_j \tilde{e}_j(x) = f(x)$

5.5.4 Connection between the measure and the embedding

Note that the measure μ on $L^2(X, \mu)$ was **arbitrary**, with the relatively strong assumption $\text{supp}(\mu) = X$. Since the kernel and the RKHS is independent of the measure, for all such measures this construction results in the same RKHS. Furthermore, if the kernel k is such that $\{e_i\}_{i \geq 1}$ is a basis of $L^2(X, \mu)$, the previous map

is an isometric isomorphism between the spaces $\mathcal{H} \leftrightarrow L^2(X, \mu)$, instead of an embedding to a subspace.

Mercer's theorem gives useful intuition into not only to the functions which make up the RKHS, but also the norm used for approximation. Suppose for example that μ is a nondegenerate probability measure. Then by definition each basis function $\tilde{e}_j \in L^2(X, \mu)$ limits an area size of the unit square. However, for a function $\sum_{j \in J} \sqrt{\lambda_j} a_j \tilde{e}_j$ in the RKHS, the functions are scaled to zero by multiplicative factors λ_j , which may decrease quite rapidly. While in principle a kernel method may use 'an infinite dimensional feature space', this guides the learning process in practise to arrive at a function determined by the amount of data and it's capacity to overwrite these factors. This observation is closely related to the effective dimension to be presented later.

Suppose this condition does not hold, $\text{supp}(\mu) \neq X$. Then for any two representers $g_1, g_2 \in [(f)]$, in the set $X_0 := X \setminus X_\mu \neq \emptyset$ it is possible that $g_1(x) \neq g_2(x)$. This motivated the association to the function restricted to $\text{supp}(\mu)$:

$$g|_{\text{supp}(\mu)} \text{ for any representer } g \in [(f)]$$

For such measure, for each equivalence class in $L^2(X, \mu)$ and the associated function $g|_{\text{supp}(\mu)}$, \mathcal{H} contains the extensions of functions $g|_{\text{supp}(\mu)}$ to X . We obtain the following hierarchy

	k universal	k \neg universal
$\text{supp}(\mu) = X$	$\mathcal{H} = L^2(X, \mu)$	$\mathcal{H} \subset L^2(X, \mu)$
$\text{supp}(\mu) \subset X$	\mathcal{H} extension $L^2(X, \mu) _{X_\mu}$	\mathcal{H} extension $\subset L^2(X, \mu) _{X_\mu}$

5.6 Convergence

The empirical operators are closely connected to the population operators. To utilize the law of large numbers in a precise way, one can use inequalities to bound the difference between the empirical mean and the expectation of random variables.

Theorem: Hoeffding inequality (Mohri;Rostamizadeh;& Ameet, 2012)

If X_1, \dots, X_n are independent real random variables such that $|X_i| \leq M$ and $\bar{X} = \frac{1}{n} \sum_i X_i$, then

$$P(|\bar{X} - E[X]| \geq \epsilon) \leq 2e^{-\frac{n\epsilon^2}{2M^2}}$$

We can reformulate this as the following bound (Abu-Mostafa;Magdon-Ismail;& Lin, 2012). With probability greater or equal to $1 - \delta$:

$$|\bar{X} - E[X]| \leq \sqrt{\frac{2M^2}{n} \log\left(\frac{2}{\delta}\right)}$$

The same inequality extends to Hilbert spaces (Rosasco;Belkin;& De Vito, 2010). If ξ_1, \dots, ξ_n are random variables in a separable real Hilbert space H such that $\|\xi_i\|_H \leq M$, then with probability greater or equal to $1 - \delta$:

$$\|\bar{X} - E[X]\|_H \leq \sqrt{\frac{2M^2}{n} \log\left(\frac{2}{\delta}\right)}$$

5.6.1 Operator and spectral converge

To obtain insight about the relation of the sample operators and population operators, consider the Hilbert space random variable $\xi(f) = f(x)K_x \in \mathcal{H}$ on (Z, ρ) . We obtain for the empirical mean of $\xi_1, \dots, \xi_m \sim \xi$ and the population mean (Smale & Zhou, 2007):

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \xi_i(f) &= \frac{1}{m} \sum_{i=1}^m f(x_i)K_{x_i} = \frac{1}{m} S_{\bar{x}}^* S_{\bar{x}} f \\ E[\xi(f)] &= \int_X K_x f(x) d\rho_X(x) = C_k f \end{aligned}$$

This shows that $C = \frac{1}{m} S_{\bar{x}}^* S_{\bar{x}}$ is an unbiased estimate of C_K , the probability of deviation decreasing with increasing number of samples as a consequence of the law of large numbers. In fact, the difference can be explicitly bounded.

Theorem (De Vito;Rosasco;Caponetto;De Giovannini;& Odone, 2005)

With probability greater or equal to $1 - \delta$:

$$\|C - C_k\|_{HS} \leq \sqrt{\frac{2\|k\|_\rho^2}{n} \log\left(\frac{2}{\delta}\right)}$$

Proof

$\xi_i(f) = \langle f, K_{x_i} \rangle K_{x_i}$ is an operator $\mathcal{H} \rightarrow \mathcal{H}$ and $\|\xi\|_{HS}^2 = \|K_x\|^4 \leq \|k\|_\rho^2$ implies $\|\xi_i\|_{HS} \leq \|k\|_\rho$. Apply Hoeffding inequality to obtain the result.

The relation between $K: \mathbb{R}^m \rightarrow \mathbb{R}^m$ and $T_k: L^2(X, \rho) \rightarrow L^2(X, \rho)$ cannot be stated through operator norm, because they do not operate on the same space. However, a simple spectral converge result can be show utilizing the following theorem referred in (Rosasco;Belkin;& De Vito, 2010).

Theorem

Let \mathcal{H} be a seperable Hilbert space with A, B self-adjoint compact operators. Let $\{\mu_i\}_{i \geq 1}$ be an enumeration of eigenvalues of $B - A$. Then there exists extended enumerations $\{\alpha_i\}_{i \geq 1}$ and $\{\beta_i\}_{i \geq 1}$ of eigenvalues of A and B such that for any nonnegative convex function ϕ with $\phi(0) = 0$:

$$\sum_{i \geq 1} \phi(\beta_i - \alpha_i) \leq \sum_{i \geq 1} \phi(\mu_i)$$

Theorem (Rosasco;Belkin;& De Vito, 2010)

If $\{\lambda_i\}_{i \geq 1}$ and $\{\mu_i\}_{i \geq 1}$ are the eigenvalues of C and C_k , they are also the eigenvalues of K and T_k , respectively, and for nonnegative convex function ϕ :

$$\sum_{i \geq 1} \phi(\mu_i - \lambda_i)^2 \leq \|C - C_k\|_{HS}^2$$

Proof

Apply the above theorem for $\phi(x) = x^2$, $A = C$, $B = C_k$.

Because $\|C - C_k\|_{HS}^2$ was bounded previously, we have bound the difference of the spectrum of K and T_k in the ℓ^2 norm.

Finally, equivalent converge results and bounds may be obtained for the operators related to the projections. For example, consider the following definitions, called the **empirical effective dimension** and the **effective dimension**:

$$D_\lambda(K) := \text{trace}[K(K + \lambda I)^{-1}] = \sum_{i \geq 1} \frac{\lambda_i}{\lambda_i + \lambda}$$

$$D_\lambda(T_k) := \text{trace}[T_k(T_k + \lambda I)^{-1}] = \sum_{i \geq 1} \frac{\mu_i}{\mu_i + \lambda}$$

Note that $D_\lambda(T_k)$, depending on the unknown measure, is a theoretical quantity, while $D_\lambda(K)$ may be computed in practise. Similar to the previous operators, converge implies we can take the empirical version as an approximation of the population version. The following result which establishes the convergence of their spectrum was proven in (Caponnetto;Rosasco;De Vito;& Verri, 2005):

For $0 < \eta < 1, \Delta > 0$ and $m \geq \left(\frac{4\|k\|_\rho}{\lambda\Delta} \left(1 + \frac{\|k\|_\rho}{\lambda}\right) \log\left(\frac{2}{\eta}\right)\right)^2$, with probability $1 - 2\eta$:

$$|D_\lambda(T_k) - D_\lambda(K)| \leq \Delta$$

5.6.2 Solution convergence

Recall the error for a hypothesis $f_{\bar{z}} = \operatorname{argmin}_{f \in \mathcal{H}} \mathcal{E}_{\bar{z}}(f)$ learned from data was defined:

$$\mathcal{E}(f_{\bar{z}}) = \mathcal{E}_{\mathcal{H}}(f_{\bar{z}}) + \mathcal{E}(f_{\mathcal{H}})$$

Where $f_{\mathcal{H}} = \operatorname{argmin}_{f \in \mathcal{H}} \mathcal{E}(f)$ and the **sampling error** $\mathcal{E}_{\mathcal{H}}(f_{\bar{z}})$ was defined:

$$\mathcal{E}_{\mathcal{H}}(f_{\bar{z}}) = \mathcal{E}(f_{\bar{z}}) - \mathcal{E}(f_{\mathcal{H}})$$

One of the most important questions in learning is whether the function $f_{\bar{z}}$ minimizing the empirical risk approaches to the function $f_{\mathcal{H}}$ minimizing the population risk in the sense that the sampling error $\mathcal{E}_{\mathcal{H}}(f_{\bar{z}}) \rightarrow 0$ as $n \rightarrow \infty$. This is called **consistency** of the algorithm.

Consistency depends on the hypothesis set, and a modern research is still refining the necessary and sufficient conditions for the consistency of empirical risk minimization. Necessary and sufficient conditions include \mathcal{H} being a Glivenko-Cantelli class of functions or having a finite V_γ dimension, and a sufficient condition for consistency is compactness of \mathcal{H} (Poggio & Smale, 2003).

Since $f_{\bar{z}}$ is a random variable, we can express consistency requirement as a **learning bound**:

$$P[\mathcal{E}(f_{\bar{z}}) - \mathcal{E}(f_{\mathcal{H}})] \leq \eta(\epsilon, n) \quad \forall \epsilon > 0, n \in \mathbb{N}$$

where $\eta(\epsilon, n)$ does not depend on ρ and $\eta(\epsilon, n) \rightarrow 0$ as $n \rightarrow \infty$.

However, to obtain such uniform estimates for the convergence, additional restrictions need to be specified on the measure ρ . The problem is that even though there always exists a solution to the empirical minimization problem which is guaranteed to converge to the target function, the convergence may be arbitrarily slow given an arbitrary probability measure ρ (Bauer;Pereverzev;& Rosasco, 2007).

5.6.3 Learning bounds for kernel RLS

Next we specialize this for the RLS algorithm where the hypothesis space is an RKHS. Recalling the factorization $\mathcal{E}(f) = \|f - f_\rho\|_\rho + \mathcal{E}(f_\rho)$, the error of a hypothesis may be studied through the norm:

$$\mathcal{E}(f_{\bar{z}}) - \mathcal{E}(f_\rho) = \|f - f_\rho\|_\rho = \|f_{\bar{z}} - f_{\mathcal{H}}\|_\rho - \|f_{\mathcal{H}} - f_\rho\|_\rho$$

where $\|f_{\bar{z}} - f_{\mathcal{H}}\|_\rho$ corresponds to the sampling error and $\|f_{\mathcal{H}} - f_\rho\|_\rho$ to the approximation error. Alternatively, since the converge in \mathcal{H} with respect to the RKHS norm $\|\cdot\|_H$ implies pointwise converge, this norm is used in some analysis (Smale & Zhou, 2007).

However, the above analysis needs caution. In general, given a RKHS \mathcal{H} the approximation error may be nonzero $\inf_{f \in \mathcal{H}} \mathcal{E}(f) > \mathcal{E}(f_\rho)$ and existence of the optimal approximator $f_{\mathcal{H}}$ is not even ensured. Denote $P: L^2(X, \rho) \rightarrow L^2(X, \rho)$ the projection into $\overline{\mathcal{R}(I_K)}$. If \mathcal{H} is dense in $L^2(X, \rho)$ then $\inf_{f \in \mathcal{H}} \mathcal{E}(f) = \mathcal{E}(f_\rho)$ and we previously saw that necessary and sufficient condition for the existence and uniqueness of the best approximation minimum norm solution is $Pf_\rho \in \mathcal{R}(I_K)$. The solution then was then defined through the pseudoinverse $f_{\mathcal{H}}^\dagger = I_K^\dagger Pf_\rho$. In this case we can replace the analysis of $\|f - f_\rho\|_\rho$ by $\|f - f_{\mathcal{H}}^\dagger\|_\rho$ (Bauer; Pereverzev; & Rosasco, 2007). However, because the pseudoinverse is not continuous it does not result in a consistent algorithm.

Consider the regularized pseudoinverse and the corresponding minimization problem using a RKHS \mathcal{H} :

$$f_{\bar{z}}^\lambda = \operatorname{argmin}_{f \in \mathcal{H}} \|S_{\bar{x}} f - \bar{y}\|_n^2 + \lambda \|f\|_H^2$$

$$f^\lambda = \operatorname{argmin}_{f \in \mathcal{H}} \|I_K f - f_\rho\|_\rho^2 + \lambda \|f\|_H^2$$

This algorithm has implicit sufficient requirements on \mathcal{H} for consistency, based on the regularization and assumptions on the kernel, such as measurability, boundedness and seperability of \mathcal{H} (Bauer; Pereverzev; & Rosasco, 2007). The conditions on ρ required for bounds are often more complicated (Caponetto & De Vito, Optimal Rates for the Regularized Least Squares Algorithm, 2007).

Different regularization schemes can be analyzed by specifying general conditions on regularization which are sufficient for the consistency of learning. Tikhonov regularization, the regularized least squares formulation, satisfies these conditions and therefore results in a consistent algorithm (Bauer; Pereverzev; & Rosasco, 2007) (De Vito; Rosasco; & Verri, Spectral Methods for Regularization in Learning Theory, 2005).

In (Cucker & Smale, 2001) it is shown that for a fixed $\lambda > 0$, regularized least squares RKHS formulation specifies a compact hypothesis set $\overline{I_K(B_R)}$ where B_R is a ball of radius R in \mathcal{H} and compactness is in $C(X)$. They also gave the following simple consistency bound derived through the Hoeffding inequality and covering numbers as a measure of an arbitrary hypothesis set \mathcal{H} :

Theorem

Suppose $|f(x) - y| \leq M$ for all $f \in \mathcal{H}$ for almost all $z \in Z$. Then

$$P_{z \in Z^m}[\mathcal{E}_{\mathcal{H}}(f_{\bar{z}}) \leq \epsilon] \leq 1 - \delta$$

where $\delta = \text{Cov\#}\left(\mathcal{H}, \frac{\epsilon}{24M}\right) e^{-\frac{m\epsilon}{288M^2}}$.

However, not only does the regularization parameter imply a nonzero approximation error $\|f_{\lambda} - f_{\rho}\|_{\rho}$, we should furthermore consider a regularization parameter λ variable in the number of samples to obtain converge to the best approximator of f_{ρ} in the whole RKHS \mathcal{H} . Both in practise and in theory, this implies an intricate balancing between the sampling and approximation errors (Poggio & Smale, 2003).

The most general learning bounds considering both the sampling error and approximation error with possibly variable λ can be very complicated. Consider the bound, where with probability $1 - \delta$:

$$\mathcal{E}(f_{\bar{z}}) - \mathcal{E}(f_{\mathcal{H}}) \leq B(\delta, \mathcal{H}, n)$$

In these bounds, the analysis focuses especially on specifying and improving the constant in the bound and the convergence rate as a function of the number of samples m . The constant and the convergence rate are both variable, depending on the algorithm, the hypothesis set and restrictions on ρ . For an example of RLS algorithm

with assumptions resulting in $O\left(\frac{1}{\sqrt{m}}\right)$ rate, see (Caponetto, Optimal Rates for Regularization Operators in Learning Theory, 2006).

The tightness of the bound depends on the properties of the hypothesis space, which has been previously quantified with VC dimension (Poggio & Smale, 2003), Rademacher averages (Steinwart & Christmann, 2008), Covering numbers (Cucker & Smale, 2001), and recently the effective dimension (Zhang, 2005). Recall we defined the effective dimension, a property of the RLS estimator using kernel $k(x, y)$, as follows:

$$D_\lambda(T_k) = \sum_i \frac{\lambda_i}{\lambda_i + \lambda}$$

where λ_i are the eigenvalues of the integral operator T_k associated to kernel k and λ is the regularization parameter, obtained as the trace of the regularized population error minimizer. Effective dimension is widely used in learning bounds in the kernel methods literature, because it allows both a simple interpretation, small constant factors and optimal convergence rates (Caponetto & De Vito, Optimal Rates for the Regularized Least Squares Algorithm, 2007). In our investigation, the simplicity is very useful for comparing different hypothesis spaces and considering modifications of a hypothesis space obtained by composition or constraints.

For a general effective dimension based learning bound, see (Caponetto & De Vito, Optimal Rates for the Regularized Least Squares Algorithm, 2007), recent advancements are discussed in (Hsu; Kakade; & Zhang, 2010)

6 Learning relations

6.1 Introduction to relations

Relation is a relationship between objects which may be considered for each pair of objects $v, v' \in \mathcal{V}$. Assume the relation is quantifiable as a function $Q(v, v')$ over the Cartesian product $(v, v') \in \mathcal{V} \times \mathcal{V}$. Then in machine learning we observe samples from the underlying relation as pairs of objects coupled with an output $\mathcal{D} = \{((v, v')_i, y_i)\}_{i=1}^n$ where $y_i = Q(v, v')$. The goal is to model the underlying relation $Q(v, v')$ based on these samples. In (Waegeman, ym., 2012) quantifiable relations are classified based on output

- $Q: \mathcal{V}^2 \rightarrow \{0,1\}$ is called a crisp relation
- $Q: \mathcal{V}^2 \rightarrow [0,1]$ graded relation, which may extended to \mathbb{R} by scaling
- $Q: \mathcal{V}^2 \rightarrow \mathbb{N}$ as an injective map is an ordinal relation, where only the order is considered

Another way to classify relations is by the properties they have. In particular (Waegeman, ym., 2012) considers the following types of relations

- $Q(v, v') = Q(v', v)$ is a symmetric relation.
- $Q(v, v') = -Q(v', v)$ is an anti-symmetric relation.
- $Q(v, v') = 1 - Q(v', v)$ for Q in $[0,1]$ is a reciprocal relation.

Since is possible to transform between anti-symmetric and reciprocal relation through a scaling operation, often one of them is taken as the focus of a paper, as in (Pahikkala;Waegeman;Tsivtsivadze;Salakoski;& De Baets, 2010).

In a kernel based relation learning algorithm we have a kernel $k: \mathcal{V}^2 \times \mathcal{V}^2 \rightarrow \mathbb{R}$ defined between pairs $v, \bar{v} \in \mathcal{V}^2$. A simple notational convenience is to write this kernel as a four argument **pairwise kernel**:

$$k((v, v'), (\bar{v}, \bar{v}')) = k(v, v', \bar{v}, \bar{v}')$$

Our discussion of an abstract kernel $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ encompass this scenario with $\mathcal{X} = \mathcal{V}^2$. However, pairwise kernels allow special considerations based on the properties of the relation or the objects being predicted. In the simplest case we are pro-

vided a kernel for the pairs or a feature mapping for a pair, which is simply the standard kernel setting. However, we may impose extra constraints defined only for relations, or build up the feature mapping for the pair from the feature mappings for the objects. Examples of relation specific constraints are the aforementioned symmetry/anti-symmetry, and a widely used technique for building a relation-predictor from object features is the Kronecker kernel. These considerations are incorporated in (Pahikkala; Waegeman; Tsivtsivadze; Salakoski; & De Baets, 2010).

6.2 Symmetric and anti-symmetric kernel

In this section we specify restrictions on a pairwise kernel with the goal of restricting the learning to only symmetric and anti-symmetric relations.

Let $\Psi(v, v')$ be a feature mapping on \mathcal{V}^2 . Consider the symmetric and anti-symmetric parts

$$\Psi(v, v') = \frac{1}{2}(\Psi(v, v') + \Psi(v', v)) + \frac{1}{2}(\Psi(v, v') - \Psi(v', v))$$

Take as new features the projections into the symmetric and anti-symmetric parts

$$\Phi_S(v, v') = \frac{1}{2}(\Psi(v, v') + \Psi(v', v))$$

$$\Phi_A(v, v') = \frac{1}{2}(\Psi(v, v') - \Psi(v', v))$$

Given a kernel $K^\Psi = \langle \Psi(v, v'), \Psi(\bar{v}, \bar{v}') \rangle$, define as new kernels the **symmetric kernel** K_S^Φ and the **anti-symmetric kernel** K_A^Φ through these features:

$$\begin{aligned} K_S^\Phi(v, v', \bar{v}, \bar{v}') &= \langle \Phi_S(v, v'), \Phi_S(\bar{v}, \bar{v}') \rangle \\ &= \frac{1}{4} \{ \langle \Psi(v, v'), \Psi(\bar{v}, \bar{v}') \rangle + \langle \Psi(v, v'), \Psi(\bar{v}', \bar{v}) \rangle + \langle \Psi(v', v), \Psi(\bar{v}, \bar{v}') \rangle \\ &\quad + \langle \Psi(v', v), \Psi(\bar{v}', \bar{v}) \rangle \} \\ &= \frac{1}{4} \{ K^\Psi(v, v', \bar{v}, \bar{v}') + K^\Psi(v, v', \bar{v}', \bar{v}) + K^\Psi(v', v, \bar{v}, \bar{v}') \\ &\quad + K^\Psi(v', v, \bar{v}', \bar{v}) \} \end{aligned}$$

$$\begin{aligned} K_A^\Phi(v, v', \bar{v}, \bar{v}') &= \langle \Phi_A(v, v'), \Phi_A(\bar{v}, \bar{v}') \rangle \\ &= \frac{1}{4} \{ \langle \Psi(v, v'), \Psi(\bar{v}, \bar{v}') \rangle - \langle \Psi(v, v'), \Psi(\bar{v}', \bar{v}) \rangle - \langle \Psi(v', v), \Psi(\bar{v}, \bar{v}') \rangle \\ &\quad + \langle \Psi(v', v), \Psi(\bar{v}', \bar{v}) \rangle \} \\ &= \frac{1}{4} \{ K^\Psi(v, v', \bar{v}, \bar{v}') - K^\Psi(v, v', \bar{v}', \bar{v}) - K^\Psi(v', v, \bar{v}, \bar{v}') \\ &\quad + K^\Psi(v', v, \bar{v}', \bar{v}) \} \end{aligned}$$

Next consider their sum, called the **permutation invariant** kernel:

$$\begin{aligned} K_{PI}(v, v', \bar{v}, \bar{v}') &= K_S(v, v', \bar{v}, \bar{v}') + K_A(v, v', \bar{v}, \bar{v}') \\ &= \frac{1}{2} \left(K^\Psi(v, v', \bar{v}, \bar{v}') + K^\Psi(v', v, \bar{v}', \bar{v}) \right) \end{aligned}$$

Even though the above kernels were defined using the feature map, this was only an intermediate step. The final kernel depended only on the given pairwise kernel K^Ψ . Given an kernel $K: \mathcal{V}^2 \times \mathcal{V}^2 \rightarrow \mathbb{R}$, these may be considered as operator projections $\Gamma_S(K) = K_S$, $\Gamma_A(K) = K_A$ and $\Gamma_{PI}(K) = K_{PI}$. In this paper these specific kernels are referred to as the symmetric, anti-symmetric and permutation invariant kernels of a given kernel K . Given an arbitrary kernel, we may define these properties as follows (Pahikkala, Viljanen, Airola, & Waegeman, 2015):

- A kernel $K_S(v, v', \bar{v}, \bar{v}')$ is a symmetric pairwise kernel if $K_S(v, v', \bar{v}, \bar{v}') = K_S(v', v, \bar{v}, \bar{v}')$
- A kernel $K_A(v, v', \bar{v}, \bar{v}')$ is a anti-symmetric pairwise kernel if $K_A(v, v', \bar{v}, \bar{v}') = -K_A(v', v, \bar{v}, \bar{v}')$
- A kernel $K_{PI}(v, v', \bar{v}, \bar{v}')$ is permutatation invariant if $K_{PI}(v, v', \bar{v}, \bar{v}') = K_{PI}(v', v, \bar{v}', \bar{v})$

6.3 Examples of pairwise kernels

The anti-symmetric kernel was originally introduced in (Pahikkala;Waegeman;Tsivtsivadze;Salakoski;& De Baets, 2010) where the authors contrasted two methods for learning pairwise relations: Kronecker kernel and transitive pairwise kernel. Using the previous definition, anti-symmetric Kronecker kernel was then formed from a given Kronecker kernel to be able to learn intransitive pairwise relations. Both kernels can be viewed in a unified theoretical framework as anti-symmetrizations of different kernels or feature mappings. Their resulting forms have shared properties which motivate our theoretical considerations of the symmetrisation/anti-symmetrization of a general abstract kernel. Our main learning theoretic result about the effects of symmetrisation/anti-symmetrisation process therefore applies to both kernels.

6.3.1 *Transitive pairwise kernel*

Given a feature representation $\phi(v)$ or a kernel K^Φ for the objects $v \in \mathcal{V}$, define a feature mapping for the relation:

$$\Psi(v, v') = \phi(v)$$

The kernel for the feature mapping is then defined as the object kernel:

$$K^\Psi(v, v', \bar{v}, \bar{v}') = K^\Phi(v, \bar{v})$$

Consider the effect of anti-symmetrization on this kernel:

$$K_A^\Psi(v, v', \bar{v}, \bar{v}') = \frac{1}{4} \{K^\Phi(v, \bar{v}) - K^\Phi(v, \bar{v}') - K^\Phi(v', \bar{v}) + K^\Phi(v', \bar{v}')\}$$

The predictor becomes:

$$h(x, x') = \langle w, \Psi(x, x') - \Psi(x', x) \rangle = \langle w, \phi(v) \rangle - \langle w, \phi(v') \rangle = f(v) - f(v')$$

This particular form for the predictor has consequences for the type of relation one can learn. In (Pahikkala; Waegeman; Tsivtsivadze; Salakoski; & De Baets, 2010) the authors recollect the following definitions and the following lemma:

Definition

A function $Q: X^2 \rightarrow [0,1]$ is called a reciprocal relation if

$$\forall (x, x') \in X^2 \quad Q(x, x') + Q(x', x) = 1$$

Definition

A reciprocal relation $Q: X^2 \rightarrow [0,1]$ is called strongly ranking representable if there exists a function $f: X \rightarrow \mathbb{R}$ and a cumulative distribution function $G: \mathbb{R} \rightarrow [0,1]$ with $G(0) = \frac{1}{2}$:

$$g(f(x), f(x')) = G(f(x) - f(x'))$$

Definition

A reciprocal relation $Q: X^2 \rightarrow [0,1]$ is called strongly stochastically transitive if

$$\begin{aligned} \forall (x_i, x_j, x_k) \in X^3 \quad & \left(Q(x_i, x_j) \geq \frac{1}{2} \wedge Q(x_j, x_k) \geq \frac{1}{2} \right) \Rightarrow Q(x_i, x_k) \\ & \geq \max(Q(x_i, x_j), Q(x_j, x_k)) \end{aligned}$$

Lemma

If $Q: X^2 \rightarrow [0,1]$ is strongly ranking representable, then it is strongly stochastically transitive.

This means that this particular kernel can only learn features which have this property, which can be a positive or negative consequence depending on the prior knowledge assumptions one should make. This was the authors motivation for presenting the following kernel.

6.3.2 *Intransitive pairwise kernel*

Given a feature mapping $\phi(x)$ for the object, define a feature mapping for the relation

$$\Psi(x, x') = \phi(x) \otimes \phi(x')$$

where \otimes denotes the Kronecker product. The kernel then becomes (Waegeman, ym., 2012):

$$\begin{aligned} K^\Psi(x_i, x'_i, x_j, x'_j) &= \langle \phi(x_i) \otimes \phi(x'_i), \phi(x_j) \otimes \phi(x'_j) \rangle \\ &= \langle \phi(x_i), \phi(x_j) \rangle \otimes \langle \phi(x'_i), \phi(x'_j) \rangle = K^\phi(x_i, x_j) K^\phi(x'_i, x'_j) \end{aligned}$$

Consider the effect of anti-symmetrization on this kernel:

$$K_A^\Psi(x_i, x'_i, x_j, x'_j) = \frac{1}{2} \left(K^\phi(x_i, x_j) K^\phi(x'_i, x'_j) - K^\phi(x'_i, x_j) K^\phi(x_i, x'_j) \right)$$

The symmetric version was introduced originally in (Ben-Hur & Stafford Noble, 2005):

$$K_S^\Psi(x_i, x'_i, x_j, x'_j) = \frac{1}{2} \left(K^\phi(x_i, x_j) K^\phi(x'_i, x'_j) + K^\phi(x'_i, x_j) K^\phi(x_i, x'_j) \right)$$

The authors (Waegeman, ym., 2012) proved that the Kronecker kernel K^Ψ can be used to learn arbitrary pairwise preference relation, given that the feature representation ϕ is powerful enough (K^ϕ universal on \mathcal{V}). Furthermore, they showed that K_A^Ψ can be used to learn arbitrary anti-symmetric pairwise preference relations. The expressibility of the hypothesis space induced by this kernel is considerable.

6.4 Approximation properties of symmetric and anti-symmetric kernels

Recalling the representer theorem, the learned function has the following expression:

$$h(v, v') = \langle w, \Phi(v, v') \rangle = \sum_{i=1}^n a_i K^\Phi(v, v', \bar{v}, \bar{v}')$$

Observing the expression of the predictor, this implies that

$$h(v', v) = h(v, v') \text{ if } K^\Phi \text{ is symmetric}$$

$$h(v', v) = -h(v, v') \text{ if } K^\Phi \text{ is anti-symmetric}$$

It is trivial to verify that for the symmetrized kernel K_S^Φ and anti-symmetrized kernel K_A^Φ

$$K_S^\Phi(v, v', \bar{v}, \bar{v}') = K_S^\Phi(v', v, \bar{v}, \bar{v}')$$

$$K_A^\Phi(v, v', \bar{v}, \bar{v}') = -K_A^\Phi(v', v, \bar{v}, \bar{v}')$$

These simple results imply the following corollary.

Corollary

The symmetric/anti-symmetric kernel can learn **only** symmetric/anti-symmetric relations.

The next important question is: what relations can they learn? This of course depends on the capabilities of the original kernel; the following definitions and approximation result give justification for their broad use.

Definition: universal kernel (Steinwart & Christmann, 2008)

A continuous kernel K on a compact metric space \mathcal{X} is called universal if the RKHS induced by K is dense in the space $\mathcal{C}(\mathcal{X})$, of all continuous functions $f: \mathcal{X} \rightarrow \mathbb{R}$. That is, for every function $f \in \mathcal{C}(\mathcal{X})$ and every $\epsilon > 0$, there exists a set of input points $\{x_i\}_{i=1}^m \in \mathcal{X}$ and real numbers $\{\alpha_i\}_{i=1}^m$, with $m \in \mathbb{N}$, such that

$$\max_{x \in V} \left\{ \left| f(x) - \sum_{i=1}^m \alpha_i K(x_i, x) \right| \right\} \leq \epsilon$$

Accordingly, the hypothesis space induced by the kernel K can approximate any function in $\mathcal{C}(\mathcal{X})$ arbitrarily well, and hence it is said to have the universal approximating property.

Definition: subset kernel

Let K be a continuous kernel K on a compact metric space \mathcal{X} and let $\mathcal{F} \subseteq \mathcal{C}(\mathcal{X})$. Denote $\mathcal{F} \subseteq \mathcal{H}(K)$ if the RKHS induced by K is dense in the subset \mathcal{F} . That is, for every function $f \in \mathcal{F}$:

$$\forall \epsilon > 0 \exists \{x_i\}_{i=1}^m \subseteq \mathcal{X}, \{\alpha_i\}_{i=1}^m \subseteq \mathbb{R}: \max_{x \in \mathcal{X}} \left\| f(x) - \sum_{i=1}^m \alpha_i K(x_i, x) \right\|$$

This means that hypothesis space induced by kernel K can approximate any function in \mathcal{F} arbitrary well.

Theorem (Pahikkala, Viljanen, Airola, & Waegeman, 2015)

Let \mathcal{F} be an arbitrary set of continuous functions and let

$$\mathcal{S} = \{t | r \in \mathcal{F}, t(v, v') = r(v, v') + r(v', v)\}$$

$$\mathcal{A} = \{t | r \in \mathcal{F}, t(v, v') = r(v, v') - r(v', v)\}$$

Be the set of symmetric and anti-symmetric functions determined by \mathcal{F} .

Let K be a kernel on \mathcal{F}^2 and K^S and K^A be the corresponding symmetric and anti-symmetric kernels. Then

$$\mathcal{F} \subseteq \mathcal{H}(K) \Rightarrow \mathcal{S} \subseteq \mathcal{H}(K^S), \mathcal{A} \subseteq \mathcal{H}(K^A)$$

Proof

For a rather long and technical proof, see (Pahikkala, Viljanen, Airola, & Waegeman, 2015).

6.5 Kernel matrices of symmetric and anti-symmetric kernels

Next we will provide the kernel matrix and integral operator forms of the predictors for the symmetrized and anti-symmetrized kernels. These kernel matrices are considered for machine learning in (Pahikkala;Airola;Stock;De Baets;& Waegeman, 2013) and they are more widely used in conjunction with Kronecker type matrix products. For an extensive discussion of their properties and use in linear algebra see the book by (Abadir & Magnus, 2005).

Definition

The **vectorizing operator** is a linear operator stacking the columns of $s \times s$ matrix into an s^2 column vector:

$$\text{vec}(M) = (M_{1,1} \ M_{2,1} \ \dots \ M_{s,1} \ M_{1,2} \ \dots \ M_{s,s})^T$$

Definition

Let A be a real $m \times n$ matrix. Define matrix P_{mn} ($mn \times mn$), called the **commutation matrix**, such that

$$P_{mn} \text{vec}(A) = \text{vec}(A^T) P_{mn}$$

When $m = n$, we write $P_{mn} = P_n$. Often the subscript is clear from the context and it is omitted entirely.

Definition

Define the **symmetrizer** S_n ($n^2 \times n^2$) and **skew-symmetrizer** A_n ($n^2 \times n^2$):

$$S_n \text{vec}(A) = \text{vec}\left(\frac{1}{2}(A + A^T)\right)$$

$$A_n \text{vec}(A) = \text{vec}\left(\frac{1}{2}(A - A^T)\right)$$

Lemma

It is easily verified that

$$S_n = \frac{1}{2}(I_{n^2} + P_n)$$

$$A_n = \frac{1}{2}(I_{n^2} - P_n)$$

The following properties also follow straightforwardly:

$$P^T = P$$

$$PP = I$$

$$SS = S \text{ and } AA = A$$

$$S^T = S \text{ and } A^T = A$$

$$SA = 0$$

These abstract properties are equivalent to the following characterization: P is a permutation matrix and S, A are orthogonal projection matrices which are orthogonal to each other.

Lemma (Abadir & Magnus, 2005)

For $M, N \in \mathbb{R}^{s \times t}$, the permutation matrix has a special application as the commuter of the Kronecker product:

$$P^{s^2}(M \otimes N) = (N \otimes M)P^{t^2}$$

Lemma (Abadir & Magnus, 2005)

The Kronecker product has the following property:

$$(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$$

Consider an abstract pairwise kernel matrix $K^\Phi(e_i, e_j) := K^\Phi(v_x, v_y, v_s, v_t)$ indexed as follows:

K^Φ		v_1			...	v_n		
		v_1	...	v_n		v_1	...	v_n
v_1	v_1							
	...							
	v_n							
...					$K^\Phi(v_x, v_y, v_s, v_t)$			
v_n	v_1							
	...							
	v_n							

Then we can address it by tensor products of the unit basis vectors:

$$K^\Phi(v_x, v_y, v_s, v_t) = (e_x^T \otimes e_y^T) K^\Phi(e_s \otimes e_t)$$

Theorem

Given a kernel matrix K^Φ indexed as above, the symmetric and anti-symmetric matrices may be expressed:

$$\begin{aligned} K_S^\Phi &= SK^\Phi S \\ K_A^\Phi &= AK^\Phi A \\ K_{PI}^\Phi &= \frac{1}{2}(SK^\Phi S + AK^\Phi A) \end{aligned}$$

Proof

For full generality let K^Φ be a $ms \times nt$ matrix

$$\begin{aligned}
A_n K^\Phi A_n(v, v', \bar{v}, \bar{v}') &= (e_v^T \otimes e_{v'}^T) A_n K^\Phi A_n(e_{\bar{v}} \otimes e_{\bar{v}'}') \\
&= \frac{1}{4} (e_v^T \otimes e_{v'}^T) (I_{n^2} - K_{ms}) K^\Phi (I_{n^2} - K_{nt}) (e_{\bar{v}} \otimes e_{\bar{v}'}') \\
&= \frac{1}{4} \{ (e_v^T \otimes e_{v'}^T) K^\Phi (e_{\bar{v}} \otimes e_{\bar{v}'}') - (e_v^T \otimes e_{v'}^T) K^\Phi K_{nt} (e_{\bar{v}} \otimes e_{\bar{v}'}') \\
&\quad - (e_v^T \otimes e_{v'}^T) K_{ms} K^\Phi (e_{\bar{v}} \otimes e_{\bar{v}'}') + (e_v^T \otimes e_{v'}^T) K_{ms} K^\Phi K_{nt} (e_{\bar{v}} \otimes e_{\bar{v}'}') \} \\
&= \frac{1}{4} \{ (e_v^T \otimes e_{v'}^T) K^\Phi (e_{\bar{v}} \otimes e_{\bar{v}'}') - (e_v^T \otimes e_{v'}^T) K^\Phi (e_{\bar{v}'}' \otimes e_{\bar{v}}) \\
&\quad - (e_{v'}^T \otimes e_v^T) K^\Phi (e_{\bar{v}} \otimes e_{\bar{v}'}') + (e_{v'}^T \otimes e_v^T) K^\Phi (e_{\bar{v}'}' \otimes e_{\bar{v}}) \} \\
&= \frac{1}{4} \{ K^\Phi(v, v', \bar{v}, \bar{v}') - K^\Phi(v, v', \bar{v}', \bar{v}) - K^\Phi(v', v, \bar{v}, \bar{v}') + K^\Phi(v', v, \bar{v}', \bar{v}) \}
\end{aligned}$$

The case for $S_n K^\Phi S_n$ is equivalent. The permutation invariant kernel follows directly from the definition.

Note that even though we used the Kronecker product to define an indexing scheme, this is unrelated to the Kronecker kernel. The kernel indexed as before is arbitrary, and the indexing was used to define the symmetric and anti-symmetric kernel matrices using the symmetrizer/anti-symmetrizer matrices applied to the original kernel matrix.

Eigensystem of the permutation invariant kernel

It is easily verified that the matrices $K_S^\Phi, K_A^\Phi, K_{PI}^\Phi$ commute. This implies that they share eigenvectors. Using diagonalization of self-adjoint matrices, they can be written:

$$\begin{aligned}
K_S^\Phi &= V \Lambda^S V^* \\
K_A^\Phi &= V \Lambda^A V^* \\
K_{PI}^\Phi &= V \Lambda^{PI} V^*
\end{aligned}$$

Lemma

In the Kronecker kernel case $K^\Phi = K^\Phi \otimes K^\Phi$, where $K^\Phi(v, v')$ contains the node evaluations.

Proof

$$\begin{aligned}
K^\Phi \otimes K^\Phi(v, v', \bar{v}, \bar{v}') &= (e_v^T \otimes e_{v'}^T) (K^\Phi \otimes K^\Phi) (e_{\bar{v}} \otimes e_{\bar{v}'}') \\
&= (e_v^T K^\Phi) \otimes (e_{v'}^T K^\Phi) (e_{\bar{v}} \otimes e_{\bar{v}'}') = (e_v^T K^\Phi e_{\bar{v}}) \otimes (e_{v'}^T K^\Phi e_{\bar{v}'}') \\
&= (e_v^T K^\Phi e_{\bar{v}}) (e_{v'}^T K^\Phi e_{\bar{v}'}') = K^\Phi(v, \bar{v}) K^\Phi(v', \bar{v}')
\end{aligned}$$

Theorem

If K^Φ is a Kronecker kernel:

$$K_S^\Phi = SK^\Phi$$

$$K_A^\Phi = AK^\Phi$$

Proof

Because $P_n(K^\Phi \otimes K^\Phi) = (K^\Phi \otimes K^\Phi)P_n$ and $P_n^2 = I$, we have:

$$\begin{aligned} A(K^\Phi \otimes K^\Phi)A &= \frac{1}{4}(K^\Phi \otimes K^\Phi - (K^\Phi \otimes K^\Phi)P_n - P_n(K^\Phi \otimes K^\Phi) + P_n(K^\Phi \otimes K^\Phi)P_n) \\ &= \frac{1}{2}(K^\Phi \otimes K^\Phi - P_n(K^\Phi \otimes K^\Phi)) = A(K^\Phi \otimes K^\Phi) \end{aligned}$$

Same applies for the symmetric matrix.

Theorem

If K^Φ is a Kronecker kernel, it is equivalent to its symmetric and anti-symmetric part, i.e. it is permutation invariant kernel:

$$K^\Phi = K_{PI}^\Phi = K_S^\Phi + K_A^\Phi$$

This makes the Kronecker kernel and its spectrum relative to the symmetric/anti-symmetric parts particularly simple.

Proof

This follows from the same property $P_n(K^\Phi \otimes K^\Phi) = (K^\Phi \otimes K^\Phi)P_n$ combined with $P_n^2 = I$.

$$K_S^\Phi + K_A^\Phi = \frac{1}{2}(K^\Phi + PK^\Phi P) = K^\Phi$$

Above we showed that if the kernel is permutation invariant, such as the Kronecker kernel is, the eigensystem is particularly simple. If kernel is not permutation invariant, there does not seem to be a straightforward statement to make about the eigensystem. This motivates the separate abstract consideration for the kernel not assumed to be permutation invariant.

6.6 Integral operators of symmetric and anti-symmetric kernels

This section uses the following definitions and theorem for the spectral analysis of the permutation invariant kernel.

Definition (Majorization)

Let $\bar{r} = (r_i)_{i=1}^{\infty}$ and $\bar{s} = (s_i)_{i=1}^{\infty}$ be sequences decreasing monotonically to 0. Then \bar{s} majorizes \bar{r} , denoted $\bar{r} < \bar{s}$, if

$$\sum_{i=1}^m r_i \leq \sum_{i=1}^m s_i \quad \forall m \in \mathbb{N} \quad \text{and} \quad \sum_{i=1}^{\infty} r_i = \sum_{i=1}^{\infty} s_i$$

For two trace class operators T_1 and T_2 on a Hilbert space with sequences of eigenvalues $\bar{\lambda}_1$ and $\bar{\lambda}_2$, denote

$$T_1 < T_2 \text{ if } \bar{\lambda}_1 < \bar{\lambda}_2$$

Definition (Doubly-stochastic operation)

Let $\mathcal{T}(\mathcal{L})$ be the Banach space of all trace class operators on a Hilbert space \mathcal{L} , i.e. space of operators complete with respect to norm $\|T\|_1 = \text{trace}(T) < \infty$. The operation $\Gamma: \mathcal{T}(\mathcal{L}) \rightarrow \mathcal{T}(\mathcal{L})$ is doubly stochastic if it is trace-preserving ($\text{trace}(T) = \text{trace}(\Gamma(T))$), unital ($I = \Gamma(I)$) and there exists a sequence $\{E_i\}_{i=1}^{\infty}$ of compact operators such that $\Gamma(T) = \sum_{i=1}^{\infty} E_i T E_i^*$.

Theorem (Uhlmann's theorem for infinite dimensional Hilbert spaces)

If T_1 and T_2 are trace-class operators on a Hilbert space, then $T_1 < T_2$ iff there exists a doubly-stochastic operation Γ such that $T_2 = \Gamma(T_1)$.

Theorem (Integral operators) (Pahikkala, Viljanen, Airola, & Waegeman, 2015)

Let K be an arbitrary pairwise kernel and K_{PI}, K_S, K_A be its permutation invariant, symmetric and anti-symmetric forms. Let $T_K, T_{K_{PI}}, T_{K_S}, T_{K_A}$ be the integral operators of the corresponding kernels. Then

$$\begin{aligned} T_{K_P} &= P^{\mu*} T_K P^{\mu} \\ T_{K_S} &= S^{\mu*} T_K S^{\mu} \\ T_{K_A} &= A^{\mu*} T_K A^{\mu} \\ T_{K_{PI}} &= \frac{1}{2} (T_K + P^{\mu*} T_K P^{\mu}) = S^{\mu*} T_K S^{\mu} + A^{\mu*} T_K A^{\mu} \end{aligned}$$

Where $P^{\mu}: L^2(\mathcal{P}^2, \mu) \rightarrow L^2(\mathcal{P}^2, \mu)$ is a permutation operator with respect to measure μ :

$$P^{\mu}: h(\bar{v}, \bar{v}') \mapsto \frac{\mu(\bar{v}', \bar{v})}{\mu(\bar{v}, \bar{v}')} h(\bar{v}', \bar{v})$$

$$P^{\mu*}: h(\bar{v}, \bar{v}') \mapsto \frac{\mu(\bar{v}, \bar{v}')}{\mu(\bar{v}', \bar{v})} h(\bar{v}', \bar{v})$$

And the symmetrizer and anti-symmetrizer with respect to the measure μ can be written:

$$\begin{aligned} S^\mu &= \frac{1}{2}(I + P^\mu) \\ A^\mu &= \frac{1}{2}(I - P^\mu) \end{aligned}$$

Proof

For the long proof see (Pahikkala, Viljanen, Airola, & Waegeman, 2015).

An important special case is when the measure is symmetric

$$\mu(\bar{v}, \bar{v}') = \mu(\bar{v}', \bar{v}) \quad \forall (\bar{v}, \bar{v}') \in \mathcal{P}^2$$

Theorem (Eigensystem analysis of symmetric μ) (Pahikkala, Viljanen, Airola, & Waegeman, 2015)

If μ is symmetric, the set of operators $\{S, A, T_{K_S}, T_{K_A}, T_{K_{PI}}\}$ commutes, which means they share an eigensystem:

$$\begin{aligned} T_{K_S} &= V\Lambda^S V^* \\ T_{K_A} &= V\Lambda^A V^* \\ T_{K_{PI}} &= V\Lambda^{PI} V^* \end{aligned}$$

In addition $T_{K_{PI}} = ST_K S + AT_K A$ and properties of S, A split the eigenvalues into two distinct blocks corresponding to symmetric and anti-symmetric functions

$$\Lambda^{PI} = \Lambda^S + \Lambda^A \text{ and } \Lambda^S \Lambda^A = 0$$

Finally, the sequence of eigenvalues of T_K majorizes the eigenvalues of $T_{K_{PI}}$:

$$T_{K_{PI}} \prec T_K$$

Proof

With symmetric μ , S and A are self-adjoint and hence orthogonal projections. Furthermore, they are orthogonal to each other:

$$SA = AS = 0$$

We can define a trace preserving and unital operation $\Gamma: \mathcal{T}(L^2(\mathcal{P}^2, \mu)) \rightarrow \mathcal{T}(L^2(\mathcal{P}^2, \mu))$ with the operators $\left\{\frac{1}{\sqrt{2}}I, \frac{1}{\sqrt{2}}P\right\}$:

$$\Gamma: T_K \mapsto \frac{1}{2}(T_K + PT_K P)$$

Majorization then follows from Uhlmann's theorem.

6.7 Effective dimension of symmetric and anti-symmetric kernels

This consideration is split into two effective dimension considerations, both of which are used in the literature. The first is defined as the trace of the finite dimensional approximator based on the kernel matrix and the second as the trace of the trace-class integral approximator based on T_K .

Theorem (finite-dimensional general case)

The effective dimension of the symmetric and anti-symmetric kernel matrices are smaller:

$$\begin{aligned} D_\lambda(SKS) &\leq D_\lambda(K) \\ D_\lambda(AKA) &\leq D_\lambda(K) \end{aligned}$$

Proof:

See **Appendix A**.

Theorem (infinite-dimensional general case) (Pahikkala, Viljanen, Airola, & Waegeman, 2015)

$$\begin{aligned} D_\lambda(T_{K_S}) &\leq D_\lambda(T_K) \\ D_\lambda(T_{K_A}) &\leq D_\lambda(T_K) \end{aligned}$$

Proof

Since A^μ is a projection matrix, $A^\mu(L^2(\mathcal{P}^2, \mu)) \subseteq L^2(\mathcal{P}^2, \mu)$, the operator $A^{\mu*}T_K A^\mu$ constrains the action of the integral operator T_K onto the range of A^μ , which is a subspace of $L^2(\mathcal{P}^2, \mu)$. The eigenfunctions ϕ_i associated with nonzero eigenvalues λ_i of T_{K_A} belong to this subspace, and are seen to satisfy:

$$T_K \phi_i - \lambda_i \phi_i = p \text{ with } p \perp A^\mu(L^2(\mathcal{P}^2, \mu))$$

Since $A^\mu(L^2(\mathcal{P}^2, \mu)) \subseteq L^2(\mathcal{P}^2, \mu)$, we can use a well known theorem (see e.g. (Aronszajn, Rayleigh-Ritz and A. Weinstein methods for approximation of eigenvalues: I. Operations in a Hilbert space, 1948) and references therein) to obtain:

$$\lambda_i^{K_A} \leq \lambda_i^K \text{ for } i = 1, 2, \dots$$

Same considerations apply to S^μ so the case with T_{K_S} goes analogously: $\lambda_i^{K^S} \leq \lambda_i^K$ for $i = 1, 2, \dots$

Each term in the effective dimension is a monotonically increasing function of λ_i , from term-wise comparison the result follows.

Theorem (Permutation-invariance with symmetric μ)

If the kernel is permutation invariant, for the kernel matrix and the integral operator with symmetric measure μ , we previously showed $\Lambda^{PI} = \Lambda^S + \Lambda^A$ and $\Lambda^S \Lambda^A = 0$. Then we have a particularly simple result for the effective dimension

$$D_\lambda(T_K) = D_\lambda(T_{K_S}) + D_\lambda(T_{K_A})$$

Theorem (Pahikkala, Viljanen, Airola, & Waegeman, 2015)

If the measure is symmetric

$$D_\lambda(T_K) \leq D_\lambda(T_{K_{PI}})$$

Proof

A following result was recently proven by Mari et al. (2014).

Let $\bar{r} = (r_i)_{i=1}^\infty$ and $\bar{s} = (s_i)_{i=1}^\infty$ be sequences decreasing monotonically to 0 with $\sum_{i=1}^\infty r_i = \sum_{i=1}^\infty s_i = 1$. Then

$$\bar{r} < \bar{s} \Leftrightarrow \sum_{i=1}^m \rho(r_i) \geq \sum_{i=1}^m \rho(s_i)$$

For all real non-negative strictly concave functions ρ defined on the segment $[0, 1]$. Since $\rho(r) = \frac{r}{r+\lambda}$ satisfies these properties, with scaling of the eigenvalues we have that

$$T_2 < T_1 \Rightarrow D_\lambda(T_1) \leq D_\lambda(T_2)$$

And the result follows from previous majorization theorem.

Summarizing the results in chapters:

- $D_\lambda(K) = D_\lambda(SKS) + D_\lambda(AKA)$ if the kernel is permutation invariant.
- $D_\lambda(T_K) = D_\lambda(T_{K_S}) + D_\lambda(T_{K_A})$ if the kernel is permutation invariant and μ is symmetric.

- $D_\lambda(T_K) \leq D_\lambda(T_{K_{PI}})$ if μ is symmetric.
- $D_\lambda(SKS) \leq D_\lambda(K)$ and $D_\lambda(AKA) \leq D_\lambda(K)$ in any case.

References

- Abadir, K.;& Magnus, J. (2005). *Matrix Algebra*. Cambridge University Press.
- Abu-Mostafa, Y.;Magdon-Ismael, M.;& Lin, H.-T. (2012). *Learning from Data - A Short Course*. AMLBook.com.
- Aronszajn. (1948). Rayleigh-Ritz and A. Weinstein methods for approximation of eigenvalues: I. Operations in a Hilbert space. *Proceedings of the National Academy of Sciences*, 474-480.
- Aronszajn. (1948). Theory of Reproducing Kernels. *AMS*, 337-404.
- Bauer, F.;Pereverzev, S.;& Rosasco, L. (2007). On regularization algorithms in learning theory. *Journal of complexity*.
- Ben-Hur, A.;& Stafford Noble, W. (2005). Kernel methods for predicting protein-protein. *Bioinformatics*, i38-i46.
- Caponetto, A. (2006). *Optimal Rates for Regularization Operators in Learning Theory*. MIT Computer Science and Artificial Intelligence Laboratory.
- Caponetto, A.;& De Vito, E. (2007). Optimal Rates for the Regularized Least Squares Algorithm. *Foundations of Computational Mathematics*, 331-368.
- Caponnetto, A.;Rosasco, L.;De Vito, E.;& Verri, A. (2005). *Empirical Effective Dimension and Optimal Rates for Regularized Least Squares Algorithm*. MIT CSAIL.
- Cucker, F.;& Smale, S. (2001). On the Mathematical Foundations of Learning. *Bulletin of the American Mathematical Society*, 1-49.
- De Vito, E.;Rosasco, L.;& Verri, A. (2005). *Spectral Methods for Regularization in Learning Theory*. Università di Genova.
- De Vito, E.;Rosasco, L.;Caponetto, A.;De Giovannini, U.;& Odone, F. (2005). Learning from Examples as an Inverse Problem. *Journal of Machine Learning Research*, 883-904.
- De Vito, E.;Rosasco, L.;Caponnetto, A.;Piana, M.;& Verri, A. (2004). Some Properties of Regularized Kernel Methods. *Journal of Machine Learning Research*, 1363-1390.

- Decoste, D.;& Bernhard, S. (2002). Training Invariant Support Vector Machines. *Machine Learning*, 161-190.
- Engl, H.;Hanke, M.;& Neubauer, A. (1996). *Regularization of Inverse Problems*. Kluwer Academic Publishers.
- Haasdonk, B.;& Keysers, D. (2002). Tangent distance kernels for Support Vector Machines. *Proceedings of the 16th International Conference on Pattern Recognition*, (ss. 864-868).
- Hein, M.;& Olivier, B. (2004). *Kernels, Associated Structures and Generalizations*. Max Planck Institute for Biological Cybernetics.
- Hsu, D.;Kakade, S.;& Zhang, T. (2010). Random Design Analysis of Ridge Regression.
- Kreyszig, E. (1989). *Introductory Functional Analysis with Applications*. John Wiley & Sons.
- Lauer, F.;& Bloch, G. (2007). Incorporating prior knowledge in support vector machines for classification - A review. *Neurocomputing*, 1578-1594.
- Mohri, M.;Rostamizadeh, A.;& Ameet, T. (2012). *Foundations of Machine Learning*. The MIT Press.
- Moon, T. (1999). *Mathematical Methods and Algorithms for Signal Processing*. Prentice Hall.
- Pahikkala, T., Viljanen, M., Airola, A., & Waegeman, W. (2015). *Spectral Analysis of Symmetric and Anti-Symmetric Pairwise Kernels*. arXiv. Retrieved from <http://arxiv.org/abs/1506.05950>
- Pahikkala, T.;Airola, A.;Stock, M.;De Baets, B.;& Waegeman, W. (2013). Efficient Regularized Least-Squares Algorithms for Conditional Ranking on Relational Data. *Machine Learning*, 321-356.
- Pahikkala, T.;Waegeman, W.;Tsvitsivadze, E.;Salakoski, T.;& De Baets, B. (2010). Learning intransitive reciprocal relations with kernel methods. *European Journal of Operational Research*, 676–685.
- Poggio, T.;& Smale, S. (2003). The Mathematics of Learning: Dealing with Data. *Notices of the AMS*.

- Poincaré, H. (1890). Sur les Equations aux Dérivées Partielles de la Physique Mathématique. *American Journal of Mathematics*, 211-294.
- Rosasco, L.;Belkin, M.;& De Vito, E. (2010). On Learning with Integral Operators. *Journal of Machine Learning Research*, 905-934.
- Schulz-Mirbach. (1994). Constructing invariant features by averaging techniques. *Proceedings of the 12th International Conference on Pattern Recognition*, (ss. 387–390). Jerusalem, Israel.
- Schölkopf, B.;& Smola, A. (2001). *Learning with Kernels*. The MIT Press.
- Smale, S.;& Zhou, D.-X. (2004). Shannon Sampling and Function reconstruction. *Bulletin of the American Mathematical Society*, 279-305.
- Smale, S.;& Zhou, D.-X. (2005). Shannon sampling II: Connections to learning theory. *Applied and Computational Harmonic Analysis*, 285-302.
- Smale, S.;& Zhou, D.-X. (2007). Learning Theory Estimates via Integral Operators and Their Approximations. *Constructive Approximation*, 153-172.
- Steinwart, I.;& Christmann, A. (2008). *Support Vector Machines*. Springer.
- Stenger, W.;& Alexander, W. (1972). *Methods of intermediate problems for Eigenvalues : theory and ramifications*. Academic Press.
- Waegeman, W.;Pahikkala, T.;Airola, A.;Salakoski, T.;Stock, M.;& Baets, B. D. (2012). A kernel-based framework for learning graded relations from data. *IEEE Transactions on Fuzzy Systems*, 1090 - 1101.
- Wolf, L.;& Shashua, A. (2003). Learning over sets using kernel principal angles. *The Journal of Machine Learning Research*, 913-931.
- Zhang, T. (2005). Learning Bounds for Kernel Regression Using Effective Data Dimensionality. *Neural Computation*, 2077-2098.

Appendix A: Effective dimension

A.1 Proof of $D_\lambda(SKS) \leq D_\lambda(K)$

Geometrical interpretation of S and A

P is a symmetric matrix with $P^2 = I$ and thus it has eigenvalues -1 and 1 . This means P has a spectral decomposition:

$$P = P_1 - P_{-1}$$

Because P_1 and P_{-1} are projection matrices to the eigenspaces of a symmetric matrix P we additionally have that

$$I = P_1 + P_{-1} \text{ and } P_1 P_{-1} = 0$$

Through this we arrive at a natural interpretation of S and A :

$$S := \frac{1}{2}(I + P) = P_1$$
$$A := \frac{1}{2}(I - P) = P_{-1}$$

This means S and A are orthogonal projections into orthogonal subspaces of V that span the space. These subspaces are spanned by eigenvectors of P associated with eigenvalues 1 and -1 , respectively.

Denoting $V_S = \text{im}(S)$ and $V_A = \text{im}(A)$, we have $\ker(S) = V_A$, $\ker(A) = V_S$ and $V_S^\perp = V_A$.

This implies $V = V_S \oplus V_A$.

Eigenvectors of SKS

Next a simple result about the eigenvectors of SKS .

Write any vector as a sum of its components in each subspace: $x = x_s + x_a$ where $x_s \in V_S, x_a \in V_A$. Denote $x' = Kx$ and write $x' = x'_s + x'_a$.

Then $SKSx = SKx_s = S(x'_s + x'_a) = x'_s$. We have $x'_s = \mu x \Leftrightarrow x'_s = \mu(x_s + x_a)$.

If $\mu \neq 0$, then $x'_s = \mu x_s$ and $x_a = 0$.

If $\mu = 0$, then $x'_s = 0$.

For eigenvectors of SKS associated with nonzero eigenvalues, the eigenvectors belong to V_S and K transforms x_s to μx_s in V_S with the component in V_A being unconstrained. For zero eigenvalues, K transforms x in V to V_A .

Nonzero eigenvalues of SKS as a principal submatrix problem

For the effective dimension we are only interested in nonzero eigenvalues. In what follows we transform this problem to another form where we can use a general theorem.

The eigenvectors of P form basis $B = \{x_{s1}, \dots, x_{sk}, x_{a(k+1)}, \dots, x_{an}\}$. Denote the standard basis as $E = \{e_1, \dots, e_n\}$.

Express $x \in V_S$ in terms of B : $x = \alpha_1 x_{s1} + \dots + \alpha_k x_{sk}$

For eigenvalues $\mu \neq 0$, we seek $\alpha_1, \dots, \alpha_k$ such that

$$Kx_s = \mu x_s + x'_a = \mu \alpha_1 x_{s1} + \dots + \mu \alpha_k x_{sk} + \beta_{k+1} x_{a(k+1)} + \dots + \beta_n x_{an}$$

We can write this in matrix form in terms of the basis B , where K_x denotes the transformation corresponding to K :

$$K_x \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_k \\ \bar{0} \end{pmatrix} = \begin{pmatrix} \mu \alpha_1 \\ \vdots \\ \mu \alpha_k \\ \bar{\beta} \end{pmatrix} \text{ coefficients w. r. t } \begin{bmatrix} x_{si} \\ x_{aj} \end{bmatrix}$$

Define $K_{(i,j,n,m)}$ as a submatrix (k_{xy}) where $x: i \dots n, y: j \dots m$. The above problem reduces to

$$\begin{pmatrix} k_{11} & \dots & k_{1k} \\ \vdots & \ddots & \vdots \\ k_{k1} & \dots & k_{kk} \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_k \end{pmatrix} = \begin{pmatrix} \mu \alpha_1 \\ \vdots \\ \mu \alpha_k \end{pmatrix} \text{ and } \begin{pmatrix} k_{(k+1)1} & \dots & k_{(k+1)k} \\ \vdots & \ddots & \vdots \\ k_{n1} & \dots & k_{nk} \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_k \end{pmatrix} = \begin{pmatrix} \beta_{k+1} \\ \vdots \\ \beta_n \end{pmatrix}$$

$$\Leftrightarrow K_{(1,1,k,k)} \bar{\alpha} = \mu \bar{\alpha} \text{ and } K_{(1,k+1,k,n)} \bar{\alpha} = \bar{\beta}$$

Since $\bar{\beta}$ is unconstrained, this is an eigenvalue and eigenvector problem of the principal submatrix:

$$K_{(1,1,k,k)} \bar{\alpha} = \mu \bar{\alpha}$$

Eigenvalues of the principal submatrix

Poincaré's separation theorem (Abadir & Magnus, 2005, s. 348):

Let A be $n \times n$ symmetric matrix with eigenvalues $\lambda_1(A) \geq \lambda_2(A) \geq \dots \geq \lambda_n(A)$ and let A_r be $r \times r$ principal submatrix of A with eigenvalues $\lambda_1(A_r) \geq \lambda_2(A_r) \geq \dots \geq \lambda_r(A_r)$. Then $\lambda_k(A_r) \leq \lambda_k(A)$ for $k = 1 \dots r$.

This theorem implies that for the eigenvalues of SKS , that is, the solutions to $K_{(1,1,k,k)} \bar{\alpha} = \mu \bar{\alpha}$ (principal submatrix of K_x): $\mu_i \leq \lambda'_i$

Define $T = (x_{s1}, \dots, x_{sk}, x_{a(k+1)}, \dots, x_{an})$, which is a basis transformation $B \rightarrow E$. Since the eigenvectors were orthogonal, we have $T^{-1} = T'$.

$K_x = TKT'$, which means that K_x is similar to K . Similar matrices share eigenvalues: $\lambda'_i = \lambda_i$ for all i .

Due to similarity, for eigenvalues λ_i of K : $\mu_i \leq \lambda_i$

The effective dimension

Each term $\frac{\lambda_i}{\lambda_i + \lambda}$ of the effective dimension is a monotonically increasing function in $\lambda_i \in [0, \infty)$

where $\lambda > 0$. Therefore $\frac{\lambda_x}{\lambda_x + \lambda} \geq \frac{\lambda_y}{\lambda_y + \lambda} \Leftrightarrow \lambda_x \geq \lambda_y$ ($\lambda_i \geq 0, \lambda > 0$). For the eigenvalues of a PSD matrix, $\lambda_i \geq 0$.

Comparing the term expansions we have

$$D_\lambda(K) = D_\lambda(K_x) = \sum_{i=1}^n \frac{\lambda_i}{\lambda_i + \lambda} = \frac{\lambda_1}{\lambda_1 + \lambda} + \dots + \frac{\lambda_r}{\lambda_r + \lambda} + \dots + \frac{\lambda_n}{\lambda_n + \lambda}$$

$$D_\lambda(SKS) = D_\lambda(K_{(1,1,k,k)}) = \sum_{i=1}^r \frac{\mu_i}{\mu_i + \lambda} = \frac{\mu_1}{\mu_1 + \lambda} + \dots + \frac{\mu_r}{\mu_r + \lambda}$$

where $\lambda_i \geq \mu_i \Rightarrow \frac{\lambda_i}{\lambda_i + \lambda} \geq \frac{\mu_i}{\mu_i + \lambda}$, $i = 1 \dots r$ and $\frac{\lambda_i}{\lambda_i + \lambda} \geq 0$, $i = (r + 1), \dots n$.

Since every term is larger, we can conclude that

$$D_\lambda(SKS) \leq D_\lambda(K)$$

A.2 Effective dimension – a simple example

This section is based on the illustrative discussion in (Zhang, 2005).

Consider the a of observations x_i , associated signals $f(x_i)$ with i.i.d distributed zero mean, σ^2 variance noise $\{n_i\}$. For each observation, we observe a corrupted response

$$y_i = f(x_i) + n_i$$

The goal is to estimate $\hat{y}_i \approx f(x_i)$ such that the mean squared error to the true signal is small:

$$\|f - \hat{y}\|_m = \frac{1}{m} \sum_{i=1}^m (f(x_i) - \hat{y}_i)^2$$

Assume now that signal vector belongs to d -dimensional space, so that for $n \times d$ projection matrix P :

$$f = Pu$$

Assume P is orthogonal: $P^T P = I_d$. The best estimator projects the observation onto the signal subspace:

$$\hat{f} = PP^T y = PP^T (Pu + n) = Pu + PP^T n$$

For this estimator, the MSE is:

$$\|\hat{f} - f\|^2 = \frac{1}{n} \|PP^T n\|^2 = \frac{1}{n} (PP^T n)^T (PP^T n) = \frac{1}{n} n^T PP^T n = \frac{1}{n} \sum_i \sum_j n_i [PP^T]_{ij} n_j$$

We know that $E[n_i] = 0$, $E[n_i^2] = \sigma^2$ and $E[n_i n_j] = E[n_i] E[n_j] = 0$:

$$E_n [\|\hat{f} - f\|^2] = \frac{\sigma^2}{n} \text{tr}(PP^T) = \frac{d}{n} \sigma^2$$

In general, u may not belong to a fixed subspace. Extend the above analysis to any linear estimator $\hat{f} = Sy$:

$$\begin{aligned} \|\hat{f} - f\|^2 &= \|S(f + n) - f\|^2 = \|Sf - f\|^2 + 2(Sf - f)^T n + n^T S^T S n \\ E_n [\|\hat{f} - f\|^2] &= \|Sf - f\|^2 + \frac{\sigma^2}{n} \text{tr}(S^T S) \end{aligned}$$

This is a bias-variance decomposition. A good estimator S optimizes the tradeoff; Sf is close to f and $\text{tr}(S^T S)$ is small.

For RLS kernel method, the estimator is

$$\hat{f} = (K + \lambda I)^{-1}Ky = Sy$$

Using the fact that (Zhang, 2005) $tr([K + \lambda I]^{-2}K^2) \leq tr([K + \lambda I]^{-1}K)$, expressing the trace as $\sum_i \frac{\lambda_i}{\lambda_i + \lambda}$, we have motivated the use of effective dimension as a measure of variance in this simple example.