

GI13/4C60: Advanced Topics in Machine Learning

Part 1: Kernels in Machine Learning

Massimiliano Pontil

1

About these lectures

Theme: We introduce the elements of the course and establish our notation. We first review a simple learning method for linear regression. This method serves to illustrate the idea of feature maps and kernels in action. We then present the notion of positive definite kernel functions, characterize some of their properties and discuss polynomial and translation invariant kernels.

Math required: Calculus, elements of linear algebra and optimization, Supervised Learning (GI01/4C55).

2

Outline

- Ridge regression
- Feature maps
- Positive definite kernels
- Kernel construction
- Kernels on Euclidean spaces

3

Notation

We summarize our notation. More symbols will be introduced during the course when needed.

$\mathcal{X}, \mathcal{Z}, \mathcal{W}$ denote sets, x, z, w their elements

\mathbb{N} is the set of natural numbers, \mathbb{R} the set of real ones

$N_n := \{1, \dots, n\}$ is the set of integers between 1 and up to n

\mathbb{R}^n is the set of n -dimensional vectors. $\mathbf{x} \in \mathbb{R}^n$ is regarded as a column vector ($n \times 1$ matrix), $\mathbf{x}' = (x_1, \dots, x_\ell) \equiv (x_i : i \in \mathbb{N}_n) \equiv (x_i)_{i=1}^n$, ($1 \times n$ matrix)

If $\mathbf{x}_1, \dots, \mathbf{x}_\ell \in \mathbb{R}^n$, \mathbf{X} denote the $\ell \times n$ matrix whose rows are the row vectors \mathbf{x}'_j . Alternatively, $\mathbf{X}' = [\mathbf{x}_1, \dots, \mathbf{x}_\ell]$

If $\mathbf{A} = (A_{ij} : i \in \mathbb{N}_n, j \in \mathbb{N}_m)$ is a $n \times m$ matrix, \mathbf{A}' denote its transpose and \mathbf{A}^+ its pseudoinverse.

\mathbf{I} (or sometimes \mathbf{I}_n) is the $n \times n$ identity matrix

If \mathcal{W} is a Hilbert space, $\langle w, z \rangle$ denotes the inner product and $\|w\| := \sqrt{\langle w, w \rangle}$ the associated norm

4

Notation (cont.)

By a kernel we mean a positive semidefinite function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

\mathbf{K} is the kernel matrix, $\mathbf{K} = (K(\mathbf{x}_i, \mathbf{x}_j) : i, j \in \mathbb{N}_\ell, \ell \in \mathbb{N}$

$\phi : \mathcal{X} \rightarrow \mathcal{W}$ is a feature map and $\{\phi(x) : x \in \mathcal{X}\}$ the feature space.

psd means *positive semidefinite*

SVD means *singular value decomposition*

5

Linear regression

Problem: We wish to find a function $g(\mathbf{x}) = \mathbf{w}'\mathbf{x}$ which best interpolates a data set $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell)\} \subseteq \mathbb{R}^n \times \mathbb{R}$

- If the data have been generated in the form $(\mathbf{x}, g(\mathbf{x}))$, the inputs are independent and $\ell = n$ then there is a unique interpolant whose parameter \mathbf{w} solves the equation

$$\mathbf{X}\mathbf{w} = \mathbf{y}$$

where $\mathbf{y} \in \mathbb{R}^\ell$ is the vector $(\mathbf{y})_i = y_i, i \in \mathbb{N}_\ell$.

- Otherwise, this problem is *ill-posed*

Remark: : If we wish to learn a non-homogeneous function $g(\mathbf{x}) = \mathbf{w}'\mathbf{x} + b$ we simply add one more dimension to \mathbf{x} and \mathbf{w} . Specifically, we replace \mathbf{x}' by the vector $(\mathbf{x}', 1)$ and \mathbf{w}' by (\mathbf{w}', b) .

6

Ill-posed problems

Informally, a problem is well-posed (in the sense of Hadamard) if

- (1) a solution exists
- (2) the solution is unique
- (3) the solution depends continuously on the data

A problem is ill-posed if it is not well-posed

Learning problems are in general ill-posed (usually because of (2)). A theme of this course is that regularization is an effective mean to make the learning problem well-posed. We will focus on regularization in a Hilbert space

7

Ridge regression

We define the square error as

$$\mathcal{L}(\mathbf{w}; S) := \sum_{i=1}^{\ell} (y_i - \mathbf{w}'\mathbf{x}_i)^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})'(\mathbf{y} - \mathbf{X}\mathbf{w})$$

and solve the variational problem

$$\min_{\mathbf{w}} \{\mathcal{L}(\mathbf{w}, S) + \lambda \mathbf{w}'\mathbf{w}\}, \quad \lambda > 0 \quad (1)$$

The parameter λ defines a trade-off between the error on the data and the norm of the vector \mathbf{w} (degree of regularization)

Setting the derivative of \mathcal{L}_λ wrt. \mathbf{w} to zero,

$$-2\mathbf{X}'(\mathbf{y} - \mathbf{X}\mathbf{w}) + 2\lambda\mathbf{w} = 0 \quad (2)$$

we obtain the *regularized solution*

$$\mathbf{w} = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_n)^{-1}\mathbf{X}'\mathbf{y} \quad (3)$$

or

$$g(\mathbf{x}) = \mathbf{w}'\mathbf{x} = \mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_n)^{-1}\mathbf{x} \quad (4)$$

8

Singular valued decomposition (review)

Let $\mathbf{X}' = [\mathbf{x}_1, \dots, \mathbf{x}_\ell]$ be the an $n \times \ell$ data matrix. Singular value decomposition (SVD) establishes that

$$\mathbf{X}'\mathbf{X} = \mathbf{U}\Lambda_n\mathbf{U}' \quad \mathbf{X}\mathbf{X}' = \mathbf{V}\Lambda_\ell\mathbf{V}'$$

where \mathbf{U} and \mathbf{V} are orthogonal matrices, $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_n]$, $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_\ell]$,

$$\Lambda_n = \text{diag}(\lambda_1, \dots, \lambda_t, \mathbf{0}_{n-t}), \quad \Lambda_\ell = \text{diag}(\lambda_1, \dots, \lambda_t, \mathbf{0}_{\ell-t}),$$

$t = \text{rank}(\mathbf{X}\mathbf{X}') = \text{rank}(\mathbf{X}'\mathbf{X})$ and $\lambda_1 \geq \dots \geq \lambda_t > 0$, $t \leq \min(\ell, n)$.

Moreover, we have

$$\mathbf{X}' = \sum_{i=1}^t \sqrt{\lambda_i} \mathbf{u}_i \mathbf{v}_i' = \mathbf{U}\Sigma\mathbf{V}'$$

where Σ is the $n \times \ell$ matrix with leading diagonal entries $\sigma_j = \sqrt{\lambda_j}$

9

Generalized solution

When λ goes to zero \mathbf{w} tends to

$$\mathbf{w}_0 := (\mathbf{X}'\mathbf{X})^+ \mathbf{X}'\mathbf{y} \quad (5)$$

where $(\mathbf{X}'\mathbf{X})^+$ is the pseudoinverse of $\mathbf{X}'\mathbf{X}$. Using SVD this is given by

$$(\mathbf{X}'\mathbf{X})^+ = \sum_{i=1}^t \sigma_i^{-1} \mathbf{u}_i \mathbf{u}_i'$$

- $g_0(\mathbf{x}) = \mathbf{w}_0' \mathbf{x}$ is called the *generalized solution*. It is the function which, among those which minimize $\mathcal{L}(\mathbf{w}; S)$, has the smallest norm of its coefficients

Dual representation

We defined the kernel function $K : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ as $K(\mathbf{x}_i, \mathbf{x}) := \mathbf{x}'_i \mathbf{x}$. The regularized solution can be written as a linear combination of ℓ kernel functions centered at the data points

$$g(\mathbf{x}) = \sum_{i=1}^{\ell} \alpha_i K(\mathbf{x}_i, \mathbf{x}) \quad (6)$$

where $\alpha = (\alpha_i : i \in \mathbb{N}_n)'$ is given by

$$\alpha = (\mathbf{K} + \lambda \mathbf{I}_\ell)^{-1} \mathbf{y} \quad (7)$$

and $\mathbf{K} := (K(\mathbf{x}_i, \mathbf{x}_j) : i, j \in \mathbb{N}_n)$

- **Function representations:** we call the functional form (or representation) $g(\mathbf{x}) = \mathbf{w}'\mathbf{x}$ the *primal form* and (1) the *dual form* (or representation)

11

Dual representation (cont.)

Proof of eqs.(6),(7): We rewrite eq.(2) as

$$\mathbf{w} = \sum_{i=1}^{\ell} \alpha_i \mathbf{x}_i \quad (8)$$

where

$$\alpha_i := \frac{y_i - \mathbf{w}'\mathbf{x}_i}{\lambda} \quad (9)$$

Consequently, we have that $\mathbf{w}'\mathbf{x} = \sum_{i=1}^{\ell} \alpha_i \mathbf{x}'_i \mathbf{x} = \sum_{i=1}^{\ell} \alpha_i K(\mathbf{x}_i, \mathbf{x})$. This proves eq.(6). Plugging eq.(8) in eq.(9) we obtain

$$\sum_{j=1}^{\ell} (K(\mathbf{x}_i, \mathbf{x}_j) + \lambda \delta_{ij}) \alpha_j = y_i, \quad \text{or, in matrix notation: } (\mathbf{K} + \lambda \mathbf{I}_\ell) \alpha = \mathbf{y}$$

from which eq.(7) follows.

Remark: in the limit of λ going to 0 α goes to $\alpha_0 = \mathbf{K}^+ \mathbf{y}$

12

Computational considerations

Training time:

- Solving for \mathbf{w} (eq.(3)) requires $O(n^3)$ operations while solving for α (eq.(7)) requires $O(n\ell^2 + \ell^3)$ operations

If $\ell \ll n$ it is more efficient to use the dual representation

Running (testing) time:

- Computing g in the primal form (eq.(6)) requires $O(n)$ operations, while the dual form (eq.(4)) requires $O(\ell n)$ operations

13

Sparse representations

Suppose each input $\mathbf{x} \in \mathbb{R}^n$ has most of its components equal to zero (eg, we consider a dataset of images where most pixels are 'black' or a dataset of short text documents represented as 'bag of words',...)

- If k denotes the number of nonzero components of the input then computing $\mathbf{x}'\mathbf{t}$ requires at most $O(k)$ operations
- If $k\ell \ll n$ (so, also $\ell, k \ll n$) the dual representation is advantageous versus the primal representation both for training ($O(k\ell^2 + \ell^3)$ vs. $O(n^3)$) and testing ($O(\ell k)$ vs. $O(n)$)

14

Feature map

The above ideas can naturally be generalized to nonlinear function regression.

By a *feature map* we mean a function $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^N$,

$$\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_N(\mathbf{x})), \quad \mathbf{x} \in \mathbb{R}^n$$

(typically $N \gg n$). $\phi(\mathbf{x})$ is called the *feature vector* and the space $\{\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^n\}$ the *feature space*

The non-linear regression function has the primal representation

$$g(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle := \sum_{j=1}^N w_j \phi_j(\mathbf{x})$$

Remark: Later we will discuss feature maps whose range is a Hilbert space (informally, an infinite dimensional space with the structure of \mathbb{R}^N)

15

Computational considerations

Again, if $\ell \ll N$ it is more efficient to work with the dual representation.

Key observation: in the dual representation we don't need to know ϕ explicitly; we just need to know the inner product between any pair of feature vectors!

Example: $N = n^2$, $\phi(\mathbf{x}) = (x_i x_j : i, j \in \mathbb{N}_n)$. In this case we have $\langle \phi(\mathbf{x}), \phi(\mathbf{t}) \rangle = (\mathbf{x}'\mathbf{t})^2$ which requires only $O(n)$ computations whereas $\Phi(\mathbf{x})$ requires $O(n^2)$ computations.

16

Kernel vs. feature map

Given a feature map ϕ we define its associated kernel function $K : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ as

$$K(\mathbf{x}, \mathbf{t}) = \langle \phi(\mathbf{x}), \phi(\mathbf{t}) \rangle, \quad \mathbf{x}, \mathbf{t} \in \mathbb{R}^n$$

- Maybe for some feature map ϕ computing $K(\mathbf{x}, \mathbf{t})$ (but not $\phi(\mathbf{x}), \phi(\mathbf{t})$) is independent on N , only depends on n

Example (cont.) If $\phi(\mathbf{x}) = (x_{i_1} x_{i_2} \cdots x_{i_d} : i_1, i_2, \dots, i_d = 1, \dots, n)$ then we have that

$$K(\mathbf{x}, \mathbf{t}) = (\mathbf{x}'\mathbf{t})^d.$$

In this case $K(\mathbf{x}, \mathbf{t})$ is computed with $O(n)$ operations, which is independent of d (and essentially of N). On the other hand, computing $\phi(\mathbf{x})$ requires $O(N)$ operations.

17

Other learning methods

Many other learning methods can be 'kernelized'. Like ridge regression they all require solving a problem where the data appears in the form of inner products and whose solution can be expressed as a linear combination of the kernel centered at the data, eg $g(\mathbf{x}) = \sum_{i=1}^{\ell} c_i K(\mathbf{x}_i, \mathbf{x})$. In particular this includes:

- Other regression or classification methods such as support vector machines, logistic regression, perceptron algorithm, Fisher discriminant analysis, etc.
- Unsupervised learning methods based on projections in feature space such as principal component analysis, partial least squares; k -means clustering, anomaly detection, etc.

We will discuss some of these methods later in the course.

18

Regularization-based learning algorithms

Let us open a short parenthesis and show that the dual form of ridge regression holds true for other loss functions. Let

$$\mathcal{L}_\lambda(\mathbf{w}; S) = \sum_{i=1}^{\ell} V(y_i, \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle) + \lambda \langle \mathbf{w}, \mathbf{w} \rangle, \quad \lambda > 0 \quad (10)$$

where $V : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is a loss function (before $(V(a, b) = (a - b)^2)$).

Theorem: If V is differentiable wrt. its second argument and \mathbf{w} is a minimizer of \mathcal{L}_λ then it has the form

$$\mathbf{w} = \sum_{i=1}^{\ell} \alpha_i \phi(\mathbf{x}_i) \quad \Rightarrow \quad g(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle = \sum_{i=1}^{\ell} \alpha_i K(x_i, x)$$

This result is usually called the *Representer Theorem*

19

Representer theorem

Setting the derivative of \mathcal{L}_λ wrt. \mathbf{w} to zero we have

$$-\sum_{i=1}^{\ell} V'(y_i, \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle) \phi(\mathbf{x}_i) + 2\lambda \mathbf{w} = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^{\ell} \alpha_i \phi(\mathbf{x}_i) \quad (11)$$

where V' is the partial derivative of V wrt. its second argument and we defined

$$\alpha_i = \frac{1}{2\lambda} V'(y_i, \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle) \quad (12)$$

Thus we conclude that

$$g(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle = \sum_{i=1}^{\ell} \alpha_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle = \sum_{i=1}^{\ell} \alpha_i K(\mathbf{x}, \mathbf{x}_i),$$

20

Some remarks

- Plugging eq.(11) in the rhs. of eq.(12) we obtain a set of equations for the coefficients α_i :

$$\alpha_i = \frac{1}{2\lambda} V' \left(y_i, \sum_{j=1}^{\ell} K(\mathbf{x}_i, \mathbf{x}_j) \alpha_j \right), \quad i = 1, \dots, \ell$$

when V is the square loss we retrieve the linear eq.(7)

- Substituting the same equation in eq.(10) we obtain an objective function for the α 's:

$$\sum_{i=1}^{\ell} V(y_i, (\mathbf{K}\alpha)_i) + \lambda \alpha' \mathbf{K} \alpha$$

Remark: the Representer Theorem holds true under more general conditions on V (for example V can be any convex function.)

21

General feature map

Let \mathcal{W} be a Hilbert space and \mathcal{X} a set. (**Remark:** we now denote by x, t, \dots the elements of \mathcal{X} . Careful! depending on the context, x_i is an element of \mathcal{X} or a component of vector $\mathbf{x} \in \mathbb{R}^n$)

The above ideas can be generalized to

- functions $g : \mathcal{X} \rightarrow \mathbb{R}$
- a feature map $\phi : \mathcal{X} \rightarrow \mathcal{W}$ and the associated kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ given by

$$K(x, t) = \langle \phi(x), \phi(t) \rangle, \quad x, t \in \mathcal{X} \quad (13)$$

In the previous slides, $\mathcal{X} = \mathbb{R}^n$ and $\mathcal{W} = \mathbb{R}^N$.

\mathcal{W} can be any Hilbert space; a frequent choice is either $\mathcal{W} = \mathbb{R}^N$ or the space of square summable sequences,

$$\mathcal{W} = \ell_2 = \{x = (x_i : i \in \mathbb{N}) : \sum_{i=1}^{\infty} x_i^2 < \infty\}$$

This is a separable Hilbert space with the inner product $\langle x, t \rangle := \sum_{i=1}^{\infty} x_i t_i$

22

Redundancy of the feature map

Warning: The feature map is not unique! For ex., if $\mathcal{W} = \mathbb{R}^N$ and ϕ generates K so does $\tilde{\phi} = U\phi$ where U is an (any!) $N \times N$ orthogonal matrix. Even the dimension of ϕ is not unique!

Example: If $n = 2$, $K(\mathbf{x}, \mathbf{t}) = (\mathbf{x}'\mathbf{t})^2$ is generated by both $\phi(\mathbf{x}) = (x_1^2, x_2^2, x_1x_2, x_2x_1)$ and $\tilde{\phi}(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$.

To emphasize this fact, instead of eq.(13) we sometime write

$$K(x, t) = \langle K_x, K_t \rangle \quad (14)$$

where $K_x : \mathcal{X} \rightarrow \mathcal{W}$ is a (anyone we like) feature map associated to K .

Remark: The quantity K_x is, at this point, only 'symbolic' but we will see later in the course that $K_x \equiv K(x, \cdot)$ (a function!) and \mathcal{W} is the reproducing kernel Hilbert space of kernel K .

23

Change of perspective

- Let us start directly with a kernel K and see when K can be expressed as an inner product in some feature space (eq.(13))

Question: Given a function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, which properties of K guarantee that there exist a Hilbert space \mathcal{W} and a feature map $\phi : \mathcal{X} \rightarrow \mathcal{W}$ such that $K(x, t) = \langle \phi(x), \phi(t) \rangle$?

24

Positive definite kernel

Definition: A function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is **positive semidefinite** if it is symmetric and the matrix $(K(x_i, x_j) : i, j \in \mathbb{N}_\ell)$ is positive semidefinite for every $\ell \in \mathbb{N}$ and every $x_1, \dots, x_\ell \in \mathcal{X}$

Remark: Some authors use the notation 'positive definite' to denote what we have called 'positive semidefinite'

Theorem: K is positive semidefinite if and only if

$$K(x, t) = \langle \phi(x), \phi(t) \rangle, \quad x, t \in \mathcal{X}$$

for some feature map $\phi : \mathcal{X} \rightarrow \mathcal{W}$

25

Positive definite kernel (cont.)

Proof of " \Leftarrow ": If $K(x, t) = \langle \phi(x), \phi(t) \rangle$ then we have that

$$\sum_{i,j=1}^{\ell} c_i c_j K(x_i, x_j) = \left\langle \sum_{i=1}^{\ell} c_i \phi(x_i), \sum_{j=1}^{\ell} c_j \phi(x_j) \right\rangle = \left\| \sum_{i=1}^{\ell} c_i \phi(x_i) \right\|^2 \geq 0$$

for every choice of $\ell \in \mathbb{N}$, $x_i \in \mathcal{X}$ and $c_i \in \mathbb{R}$, $i \in \mathbb{N}_\ell$.

The proof of ' \Rightarrow ' will be given later in the course where we will also discuss reproducing kernel Hilbert spaces.

26

Kernel construction

Which operations/combinations (eg, products, sums, composition, etc.) of a given set of kernels is still a kernel?

If we address this question we can build more interesting kernels starting from simple ones.

Example: We have already seen that $K(\mathbf{x}, \mathbf{t}) = (\mathbf{x}'\mathbf{t})^d$ is a kernel. For which class of functions $p : \mathbb{R} \rightarrow \mathbb{R}$ is $p(\mathbf{x}'\mathbf{t})$ a kernel? More generally, if K is a kernel when is $p(K(\mathbf{x}, \mathbf{t}))$ a kernel?

27

General linear kernel

If \mathbf{A} is a $n \times n$ psd matrix the function $K : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ defined by

$$K(\mathbf{x}, \mathbf{t}) = \mathbf{x}'\mathbf{A}\mathbf{t}$$

is a kernel

Proof: Since \mathbf{A} is psd we can write it in the form $\mathbf{A} = \mathbf{R}\mathbf{R}'$ for some $n \times n$ matrix \mathbf{R} . Thus K is represented by the feature map $\phi(\mathbf{x}) = \mathbf{R}'\mathbf{x}$.

Alternatively, note that:

$$\sum_{ij} c_i c_j \mathbf{x}'_i \mathbf{A} \mathbf{x}_j = \sum_{ij} c_i c_j (\mathbf{R}'\mathbf{x}_i)' (\mathbf{R}'\mathbf{x}_j) = \left\| \sum_i c_i \mathbf{R}'\mathbf{x}_i \right\|^2 \geq 0$$

28

Kernel composition

More generally, if $K : \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}$ is a kernel and $\phi : \mathcal{X} \rightarrow \mathbb{R}^N$, then

$$\tilde{K}(x, t) = K(\phi(x), \phi(t)), \quad x, t \in \mathcal{X}$$

is a kernel on \mathcal{X} .

Proof: By hypothesis K is a kernel and, so, for every x_1, \dots, x_ℓ the matrix $(K(\phi(x_i), \phi(x_j))) : i, j \in \mathbb{N}_\ell$ is psd

In particular, the above example corresponds to $K(\mathbf{x}, \mathbf{t}) = \mathbf{x}'\mathbf{t}$ and $\phi(\mathbf{x}) = \mathbf{R}'\mathbf{x}$

29

Kernel construction (cont.)

Question: If K_1, \dots, K_q are kernels on \mathcal{X} and $F : \mathbb{R}^q \rightarrow \mathbb{R}$, when is the function

$$F(K_1(x, t), \dots, K_q(x, t)), \quad x, t \in \mathcal{X}$$

a kernel?

Equivalently: when for every choice of $\ell \in \mathbb{N}$ and $\mathbf{A}_1, \dots, \mathbf{A}_q$ $\ell \times \ell$ psd matrices, is the following matrix psd?

$$(F(A_{1,ij}, \dots, A_{q,ij}) : i, j = 1, \dots, \ell)$$

We discuss some examples of functions F for which the answer to these questions is YES.

30

Convex combination of kernels

If $\lambda_j \geq 0$, $j = 1, \dots, q$ then $\sum_{j=1}^q \lambda_j K_j$ is a kernel.

This fact is immediate (a non-negative combination of psd matrices is still psd)

Example: Let $q = n$, $\mathcal{X} = \mathbb{R}^n$ and $K_i(x, t) = x_i t_i$.

In particular, this implies that

- aK_1 is kernels if $a \geq 0$
- $K_1 + K_2$ is kernels

31

Product of kernels

The pointwise product of two kernels K_1 and K_2

$$K(x, t) := K_1(x, t)K_2(x, t), \quad x, t \in \mathcal{X}$$

is a kernel.

Proof: We need to show that if \mathbf{A} and \mathbf{B} are psd matrices, so is $\mathbf{C} = (A_{ij}B_{ij} : i, j \in \mathbb{N}_\ell)$ (\mathbf{C} is also called the Schur product of \mathbf{A} and \mathbf{B}). We write \mathbf{A} and \mathbf{B} in their singular value form, $\mathbf{A} = \mathbf{U}\Sigma\mathbf{U}'$, $\mathbf{B} = \mathbf{V}\Lambda\mathbf{V}'$ where \mathbf{U}, \mathbf{V} are orthogonal matrices and $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_\ell)$, $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_\ell)$. We have

$$\begin{aligned} \sum_{i,j=1}^{\ell} a_i a_j C_{ij} &= \sum_{ij} a_i a_j \sum_r \sigma_r U_{ir} U_{jr} \sum_s \lambda_s V_{is} V_{js} \\ &= \sum_{rs} \sigma_r \lambda_s \sum_i a_i U_{ir} V_{is} \sum_j a_j U_{jr} V_{js} \\ &= \sum_{rs} \sigma_r \lambda_s \left(\sum_i a_i U_{ir} V_{is} \right)^2 \geq 0 \end{aligned}$$

32

Summary of kernel properties

The above results can be summarized as follows:

If K_1, K_2 are kernels on a set \mathcal{X} , $a \geq 0$, K a kernel on \mathbb{R}^N and $\phi : \mathcal{X} \rightarrow \mathbb{R}^N$ then the following functions are kernels on \mathcal{X}

1. $K_1(x, t) + K_2(x, t)$
2. $aK_1(x, t)$
3. $K_1(x, t)K_2(x, t)$
4. $K(\phi(x), \phi(t))$

33

Polynomial of kernels

Let $F = p$ where $p : \mathbb{R}^q \rightarrow \mathbb{R}$ is a polynomial with nonnegative coefficients. By property 1 and 2 above we conclude that p is a valid function.

In particular if $q = 1$,

$$\sum_{i=1}^d a_i (K(x, t))^i$$

is kernel if $a_1, \dots, a_d \geq 0$

34

Polynomial kernels

We now discuss kernels on \mathbb{R}^n . The above observation implies that if $p : \mathbb{R} \rightarrow \mathbb{R}$ is a univariate polynomial with nonnegative coefficients then $p(\mathbf{x}'\mathbf{t}), \mathbf{x}, \mathbf{t} \in \mathbb{R}^n$ is kernel on \mathbb{R}^n . In particular if $a \geq 0$ the following are valid polynomial kernels

- $(\mathbf{x}'\mathbf{t})^d$
- $(a + \mathbf{x}'\mathbf{t})^d$
- $\sum_{i=0}^d \frac{a^i}{i!} (\mathbf{x}'\mathbf{t})^i$

35

'Infinite polynomial' kernel

If in the last equation we set $d = \infty$ the series

$$\sum_{i=0}^{\infty} \frac{a^i}{i!} (\mathbf{x}'\mathbf{t})^i$$

converges everywhere uniform to $\exp(a\mathbf{x}'\mathbf{t})$ showing that this function is also a kernel.

Assume for simplicity that $n = 1$. A feature map corresponding to the kernel $\exp(ax\mathbf{t})$ is

$$\phi(x) = \left(1, \sqrt{ax}, \sqrt{\frac{a}{2}}x^2, \sqrt{\frac{a^3}{6}}x^3, \dots \right) = \left(\sqrt{\frac{a^i}{i!}}x^i : i \in \mathbb{N} \right)$$

- The feature space has an infinite number of dimensions!

36

Translation invariant and radial kernels

We say that a kernel $K : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is

- *Translation invariant* if it has the form

$$K(\mathbf{x}, \mathbf{t}) = H(\mathbf{x} - \mathbf{t}), \quad \mathbf{x}, \mathbf{t} \in \mathbb{R}^n$$

where $H : \mathbb{R}^n \rightarrow \mathbb{R}$ is a differentiable function.

- *Radial* if it has the form

$$K(\mathbf{x}, \mathbf{t}) = h(\|\mathbf{x} - \mathbf{t}\|), \quad \mathbf{x}, \mathbf{t} \in \mathbb{R}^n$$

where $h : [0, \infty) \rightarrow [0, \infty)$ is a differentiable function

37

The Gaussian kernel

An important example of radial kernel is the Gaussian

$$K(\mathbf{x}, \mathbf{t}) = \exp(-\beta\|\mathbf{x} - \mathbf{t}\|^2), \quad \beta > 0, \mathbf{x}, \mathbf{t} \in \mathbb{R}^n$$

It is a kernel because it is the product of two kernels

$$K(\mathbf{x}, \mathbf{t}) = \left(\exp(-\beta(\mathbf{x}'\mathbf{x} + \mathbf{t}'\mathbf{t})) \right) \exp(2\beta\mathbf{x}'\mathbf{t})$$

(We saw before that $\exp(2\beta\mathbf{x}'\mathbf{t})$ is a kernel. Clearly $\exp(-\beta(\mathbf{x}'\mathbf{x} + \mathbf{t}'\mathbf{t}))$ is a kernel with one dimensional feature map $\phi(\mathbf{x}) = \exp(-\beta\mathbf{x}'\mathbf{x})$ is indeed a kernel.)

Exercise: Can you find a feature map representation for the Gaussian kernel?

38

Periodic kernels

These are a special case of translation invariant kernels

Take $\mathcal{X} = \mathbb{R}$ and $K(x, y) = H(x - y)$, where $H : \mathbb{R} \rightarrow \mathbb{R}$ is differentiable, even and 2π -periodic. Since K is symmetric, H is even ($H(x) = H(-x)$) and its Fourier series consists of cosine's only:

$$H(x) = \sum_{n=0}^{\infty} a_n \cos(nx).$$

Then we have

$$K(x, y) = H(x - y) = a_0 + \sum_{n=1}^{\infty} a_n \sin(nx) \sin(ny) + \sum_{n=1}^{\infty} a_n \cos(nx) \cos(ny)$$

which, assuming $a_n \geq 0$, is of the form $\langle \phi(x), \phi(y) \rangle$ with

$$\phi(x) \equiv (\sqrt{a_0}, \sqrt{a_1} \sin(x), \sqrt{a_1} \cos(x), \sqrt{a_2} \sin(2x), \sqrt{a_2} \cos(2x), \dots)$$

Remark: Again the feature space is infinite dimensional. We will see that if $g(x) = \langle \mathbf{w}, \phi(x) \rangle$, $\|\mathbf{w}\|$ measures the smoothness of the function...

39

Bibliography

These lectures are based on Chapter 2 and 3 of:

- J. Shawe-Taylor and N. Cristianini. Kernel Methods for Pattern Analysis. Cambridge University Press, 2004.

40