# Kernel functions    (Quick recap)

## Feature Mapping

$$\tilde{x} = \phi(x)$$

vector in the transformed space e.g., 3-D

vector function

vector in the original space e.g., 2-D

## Kernel function:

scalar

$$k(x, x') \triangleq \phi^T(x)\phi(x')$$

two vectors in the original space

first   second (original space)

$$\phi^T(x) \quad \phi(x') \quad = \quad \phi^T(x') \quad \phi(x)$$

symmetric

## Some common examples of Kernels

**1) linear kernel**    $\phi(x) = x$ itself.

e.g., linear decision boundary $\underline{w}^T x + b = 0$

linearity in the original space

In the transformed domain

$$\tilde{w}^T \phi(x) + b = 0$$

appropriate dimension

linearity in the transformed space.

$\phi_3(x)$, $\phi_2(x)$, $\phi_1(x)$

(for a linear kernel: the transformation is an identity transformation)

② Stationary Kernel $k(\underline{x}, \underline{x}') = k(\underline{x} - \underline{x}')$

invariant to translations in the input pattern space.

e.g., music information retrieval male, female voice ('Pitch transposition')

②a RBF $k(\underline{x}, \underline{x}') = k(\|\underline{x} - \underline{x}'\|)$

(Radial Basic Function)

Key point: 'Kernel Trick'

(*) Kernel function: Why? Feature transform-decision boundary: linear in a higher dimension, or at least the linear decision boundary in the transformed space could give a better separation as compared to the original space.

(*) Compution: 'trick': make computations in the lower dimensional space itself.

# Dual Representations (Regression)

(*) Many linear models for classification and regression: which can be reformulated in terms of a dual representation in which kernels arise naturally.

## Regularised Linear Regression

$$J(\underline{w}) = \frac{1}{2} \sum_{i=1}^{N} \{ \underline{w}^T \underline{\phi}(\underline{x}) - t_i \}^2 + \frac{\lambda}{2} \underline{w}^T \underline{w}$$

"cosmetic purposes"

"cosmetic purposes"

function, to minimise

in $\phi(\underline{x})$ space

model (linear)

target value

difference b/w the target and the model

"fidelity" term

Summation: for all training data points

Regulariser

+/- discrepancies b/w the model and the target are treated in the same way

Lagrange multiplier

Regulariser → drive the system to favour low weights, unless it is supported by the data (fidelity)

Recap $\dfrac{\partial (\underline{x}^T \underline{a})}{\partial \underline{x}} = \underline{a}$

$\equiv \dfrac{\partial (\underline{a}^T \underline{x})}{\partial \underline{x}} = \underline{a}$

optimum ⟶ minimum

$$\frac{\partial J(\underline{w})}{\partial \underline{w}} = 0$$

$\lambda \geq 0$
Lagrange Multiplier

$$\frac{\partial J(\underline{w})}{\partial \underline{w}} = 0 \Rightarrow \frac{1}{2} \times \cancel{2} \sum_{i=1}^{N} \{\underline{w}^T \underline{\phi}(\underline{x}_i) - t_i\} \underline{\phi}(\underline{x}_i)$$

$$+ \frac{\cancel{2}}{\cancel{2}} \cdot \cancel{\lambda} \, \underline{w} = 0$$

$$\Rightarrow \underline{w} = \boxed{\left(-\frac{1}{\lambda}\right)} \sum_{i=1}^{N} \{\underline{w}^T \underline{\phi}(\underline{x}_i) - t_i\} \underline{\phi}(\underline{x}_i)$$

$$\underbrace{\qquad\qquad\qquad}_{a_i}$$

$$a_i \triangleq \left(-\frac{1}{\lambda}\right) \{\underline{w}^T \underline{\phi}(\underline{x}_i) - t_i\}$$

$$\Rightarrow \underline{w} = \sum_{i=1}^{N} a_i \, \underline{\phi}(\underline{x}_i)$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{\text{inner product representation}}$$

$$[\underline{\phi}(\underline{x}_1) \; \underline{\phi}(\underline{x}_2) \; \cdots \; \underline{\phi}(\underline{x}_N)] \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{bmatrix}$$



$$\underbrace{\qquad}_{\underline{a}}$$

$\Phi$ : transpose of the "$\underline{\underline{\Phi}}$"  'Design Matrix', $\underline{\underline{\Phi}}$

$\Phi$ was $\overset{N}{\updownarrow} \begin{bmatrix} \underline{\phi}^T(\underline{x}_1) \\ \vdots \\ \underline{\phi}^T(\underline{x}_N) \end{bmatrix}$ $\overset{\leftarrow M \rightarrow}{\boxed{\phantom{xx}}}$ $\overset{\leftarrow M \rightarrow}{\boxed{\phantom{xx}}}$   $\overset{\leftarrow M \rightarrow}{N \updownarrow \boxed{\Phi}}$

$$\Rightarrow \boxed{\underline{w} = \Phi^T \underline{a}}$$

Reformulate the problem in terms of $\underline{a}$, instead of $\underline{w}$

Substitute $\underline{w} = \Phi^T \underline{a}$ into the expression for $J(\underline{w})$ to try to eliminate $w$ altogether, and attempt to replace $J(\underline{w})$ with an expression which involves $\underline{a}$ alone, at the optimum.

$$J(\underline{a}) = \frac{1}{2}\sum_{i=1}^{N}\left\{ (\Phi^T\underline{a})^T \underline{\phi}(\underline{x}_i) - t_i \right\}^2 + \frac{\lambda}{2}(\Phi^T\underline{a})^T(\Phi^T\underline{a})$$

$$= \frac{1}{2}\sum_{i=1}^{N}\left\{ (\underline{a}^T\Phi)\underline{\phi}(\underline{x}_i) \right\}^2 - \frac{1}{2}2\sum_{i=1}^{N}(\underline{a}^T\Phi)\underline{\phi}(\underline{x}_i)t_i$$

$$+ \frac{1}{2}\sum_{i=1}^{N}t_i^2 + \underbrace{\frac{\lambda}{2}\underline{a}^T\Phi\Phi^T\underline{a}}$$

fourth term, is in its final form (we will see this later!)

The third term $= \frac{1}{2}\sum_{i=1}^{N}t_i^2$ we write this as an inner product

$$= \frac{1}{2}\begin{bmatrix} t_1 & t_2 & \cdots & t_N \end{bmatrix}\begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix} = \frac{1}{2}\underline{t}^T\underline{t}$$

The second term $= -\sum_{i=1}^{N}\boxed{\underline{a}^T\Phi}\underline{\phi}(\underline{x}_i)t_i$

$$= -(\underline{a}^T\Phi)\sum_{i=1}^{N}t_i\underline{\phi}(\underline{x}_i)$$

$$= -(\underline{a}^T\Phi)\underbrace{\begin{bmatrix} \underline{\phi}(\underline{x}_1) & \underline{\phi}(\underline{x}_2) \cdots \underline{\phi}(\underline{x}_N)\end{bmatrix}}_{\Phi^T}\begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix} = -\underline{a}^T\Phi\Phi^T\underline{t}$$
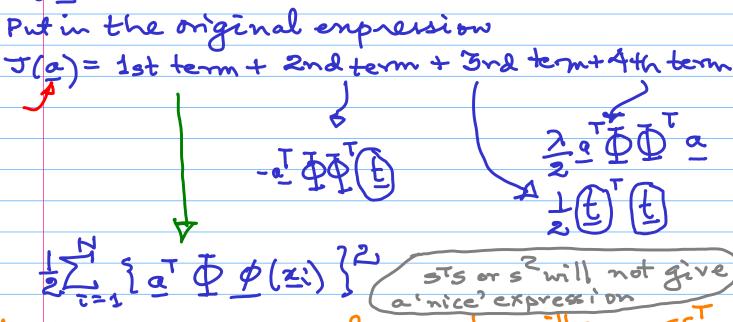
$\underline{t}$

# Basic Philosophy

(*) In some cases, formulating an original (primal) problem completely in terms of other variables (dual) is possible. There is no guarantee that given a primal problem, it should be possible to formulate a dual one in terms of another variable

(*) Even if one is able to formulate a dual problem, there is no guarantee that the dual problem may have a 'better' solution: better in terms of the computational complexity, attractiveness in terms of a kernel trick.

**Recap:**

$$J(\underline{w}) = \frac{1}{2}\sum_{i=1}^{N}\left\{\underline{w}^T\underline{\phi}(\underline{x})-t_i\right\}^2 + \frac{\lambda}{2}\underline{w}^T\underline{w}$$

$$\frac{\partial J(\underline{w})}{\partial \underline{w}} = 0 \xrightarrow{gave} \underline{w}=\Phi^T\underline{a}, \quad a_i \triangleq \frac{1}{\lambda}\left[\underline{w}^T\underline{\phi}(\underline{x}_i)-t_i\right]$$

Put in the original expression

$$J(\underline{a})= 1st\ term + 2nd\ term + 3rd\ term + 4th\ term$$

$$-\underline{a}^T\Phi\Phi^T(\underline{t})$$

$$\frac{\lambda}{2}\underline{a}^T\Phi\Phi^T\underline{a}$$

$$\frac{1}{2}(\underline{t})^T(\underline{t})$$

$$\frac{1}{2}\sum_{i=1}^{N}\left\{\underline{a}^T\Phi\,\underline{\phi}(\underline{x}_i)\right\}^2$$

$s^Ts$ or $s^2$ will not give a 'nice' expression

trick: the square of a scalar can be written as $ss^T$

$$= \frac{1}{2} \sum_{i=1}^{N} \left\{ \underline{a}^T \Phi \, \underline{\phi}(\underline{z_i}) \right\} \left\{ \underline{a}^T \Phi \, \underline{\phi}(\underline{z_i}) \right\}^T$$

$$= \frac{1}{2} \sum_{i=1}^{N} \underline{a}^T \Phi \, \underline{\phi}(\underline{z_i}) \, \underline{\phi}^T(\underline{z_i}) \, \Phi^T \underline{a}$$

Take the parts not involved in the summation, outside

$$= \frac{1}{2} \underline{a}^T \Phi \left\{ \sum_{i=1}^{N} \underline{\phi}(\underline{z_i}) \underline{\phi}^T(\underline{z_i}) \right\} \Phi^T \underline{a}$$

consider this summation alone

write as an inner product in one of the two possible ways

$$\underbrace{\left[ \underline{\phi}(\underline{z_1}) \; \underline{\phi}(\underline{z_2}) \cdots \underline{\phi}(\underline{z_N}) \right]}_{\Phi^T} \underbrace{\left[ \begin{array}{c} \underline{\phi}^T(\underline{z_1}) \\ \underline{\phi}^T(\underline{z_2}) \\ \vdots \\ \underline{\phi}^T(\underline{z_N}) \end{array} \right]}_{} \right\} \Phi$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{\Phi^T \Phi}$$

$$\Rightarrow \text{the first term} = \frac{1}{2} \underline{a}^T \left( \Phi \Phi^T \right) \left( \Phi \Phi^T \right) \underline{a}$$

The complete expression at the optimal value becomes

$$J(\underline{a}) = \frac{1}{2} \underline{a}^T \left( \Phi \Phi^T \right) \left( \Phi \Phi^T \right) \underline{a} - \underline{a}^T \left( \Phi \Phi^T \right) \underline{t}$$

$$+ \frac{1}{2} \underline{t}^T \underline{t} + \frac{\lambda}{2} \underline{a}^T \left( \Phi \Phi^T \right) \underline{a}$$

We define the **Gram Matrix** $K \triangleq \Phi \Phi^T$

what is this?

$$\Phi \Phi^T = \begin{bmatrix} \phi^T(z_1) \\ \phi^T(z_2) \\ \vdots \\ \phi^T(z_N) \end{bmatrix} \begin{bmatrix} \phi(z_1) & \phi(z_2) & \cdots & \phi(z_N) \end{bmatrix}$$

$N \times M$  $M \times N$ → $N \times N$

$M \times N$

$N \times M$

Now, what is $K(i,j)$? i'th row $\times$ j'th column

$$K(i,j) = \underbrace{\phi^T(z_i)}_{1 \times M} \underbrace{\phi(z_j)}_{M \times 1}$$

this is symmetric, scalar!

$$= k(z_i, z_j) \quad \text{the kernel function}$$

$$K(i,j) = \phi^T(z_i) \phi(z_j) = k(z_i, z_j)$$

$$J(a) = \frac{1}{2} a^T K K a - a^T K \underline{t} + \frac{1}{2} \underline{(t)}^T \underline{(t)} + \frac{\lambda}{2} a^T K a$$

[dual]

optimisation theory → $\dfrac{\partial J(a)}{\partial a} = 0$

Use result: for a quadratic form

$$\frac{\partial}{\partial a}(a^T K a) = 2 K a$$

$$\frac{\partial J(a)}{\partial a} = \frac{1}{2} \cdot 2 K K a - K \underline{(t)} + \frac{\lambda}{2} \cdot 2 K a = 0$$

$$\Rightarrow K \underline{(t)} = K(K + \lambda I_N) a$$

assume $K$ to be invertible $\quad a = (K + \lambda I_N)^{-1} \underline{(t)}$