

## 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:

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

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

Put in the original expression

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

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

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

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

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

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

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

(2)

$$\begin{aligned}
&= \frac{1}{2} \sum_{i=1}^N \{ \underline{a}^T \Phi \underline{\phi}(z_i) \} \{ \underline{a}^T \Phi \underline{\phi}(z_i) \}^T \\
&= \frac{1}{2} \sum_{i=1}^N \underline{a}^T \Phi \underbrace{\underline{\phi}(z_i) \underline{\phi}^T(z_i)} \Phi^T \underline{a}
\end{aligned}$$

Take the parts not involved in the summation, outside

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

consider this summation alone

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

$$\underbrace{[\underline{\phi}(z_1) \underline{\phi}(z_2) \dots \underline{\phi}(z_N)]}_{\Phi^T} \underbrace{\begin{bmatrix} \underline{\phi}^T(z_1) \\ \underline{\phi}^T(z_2) \\ \vdots \\ \underline{\phi}^T(z_N) \end{bmatrix}}_{\Phi} \Phi$$

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

the complete expression at the optimal value becomes

$$\begin{aligned}
J(\underline{a}) &= \frac{1}{2} \underline{a}^T (\Phi \Phi^T) (\Phi \Phi^T) \underline{a} - \underline{a}^T (\Phi \Phi^T) \underline{t} \\
&\quad + \frac{1}{2} \underline{t}^T \underline{t} + \frac{\lambda}{2} \underline{a}^T (\Phi \Phi^T) \underline{a}
\end{aligned}$$

3

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

What is this?

$$\underbrace{\Phi \Phi^T}_{N \times N} = \underbrace{\begin{bmatrix} \phi^T(z_1) \\ \phi^T(z_2) \\ \vdots \\ \phi^T(z_N) \end{bmatrix}}_{N \times M} \underbrace{[\phi(z_1) \phi(z_2) \dots \phi(z_N)]}_{M \times N}$$

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

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

this is symmetric, scalar!  
 $= k(z_i, z_j)$  the kernel function

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

$$J(\underline{a}) = \underbrace{\frac{1}{2} \underline{a}^T K K \underline{a}}_{\text{[dual]}} - \underbrace{\underline{a}^T K \underline{a}}_{\text{[dual]}} + \frac{1}{2} \underbrace{\begin{pmatrix} \underline{t} \end{pmatrix}^T \begin{pmatrix} \underline{t} \end{pmatrix}}_{\text{[dual]}} + \underbrace{\frac{\lambda}{2} \underline{a}^T K \underline{a}}_{\text{[dual]}}$$

optimisation theory  $\rightarrow \frac{\partial J(\underline{a})}{\partial \underline{a}} = 0$

Use result: for a quadratic form

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

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

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

assume  $K$  to be invertible  $\underline{a} = (K + \lambda I_N)^{-1} \begin{pmatrix} \underline{t} \end{pmatrix}$   
 [Please go to p(4)]

(4)

What is the regression?  $\underline{w}^T \phi(\underline{x}) \rightarrow$  our model  $y(\underline{x})$

$y(\underline{x}) = \underline{w}^T \phi(\underline{x})$ . For the training data, we are given target values  $t_i$ .  
 scalar ↑ input

$y(\underline{x}_i) = \underline{w}^T \phi(\underline{x}_i)$  is our modelled output for which Mother Nature (physical process) has given a value  $t_i$ .

i.e., for a 'good' model,  $y(\underline{x}_i) = \underline{w}^T \phi(\underline{x}_i)$  should be close to  $t_i$ .

$$y(\underline{x}) = \underline{w}^T \phi(\underline{x}) = (\Phi \underline{a})^T \phi(\underline{x}) = \underline{a}^T \boxed{\Phi \phi(\underline{x})}$$

$1 \times N$      $N \times N$      $N \times 1$   
 $\searrow$      $\searrow$      $\searrow$   
 scalar  $1 \times 1$

consider

$$\Phi \phi(\underline{x}) = \begin{bmatrix} \phi^T(\underline{x}_1) \\ \phi^T(\underline{x}_2) \\ \vdots \\ \phi^T(\underline{x}_N) \end{bmatrix} \quad \phi(\underline{x}) = \begin{bmatrix} \phi^T(\underline{x}_1) \phi(\underline{x}) \\ \phi^T(\underline{x}_2) \phi(\underline{x}) \\ \vdots \\ \phi^T(\underline{x}_N) \phi(\underline{x}) \end{bmatrix}$$

$N \times 1$

$$= \begin{bmatrix} k(\underline{x}_1, \underline{x}) \\ k(\underline{x}_2, \underline{x}) \\ \vdots \\ k(\underline{x}_N, \underline{x}) \end{bmatrix} \rightsquigarrow \underline{k}(\underline{x}) \quad y(\underline{x}) = \underline{a}^T \underline{k}(\underline{x}) = \underline{k}^T(\underline{x}) \underline{a}$$

$$\Rightarrow y(\underline{x}) = \underline{k}^T(\underline{x}) (\underline{K} + \lambda I_N)^{-1} \underline{t}$$

[please go]

## Physical Significance :

(\*) The dual formulation allows us to express the solution entirely in terms of the kernel function

(\*) We recover the original formulation for  $\underline{w}$ : the solution for  $\underline{a}$  can be expressed as a linear combination of the elements of  $\Phi(\underline{x})$

(\*) The prediction at  $\underline{x}$  is a linear combination of the target values from the training set.

(\*) complexity of the primal and dual formulations

$$\text{primal: } \underline{w} = \underbrace{\Phi^T}_{N \times 1} \underline{a}$$

$N \times 1$        $N \times N$        $N \times 1$   
 $N \times 1$

Solving for  $\underline{w}$  will typically involve inverting an  $N \times N$  matrix, and

typically,  $N \gg n$

$$\text{Dual: } \underline{a} = (K + \lambda I_N)^{-1} \underline{t}$$

$N \times N$  matrix

→ complexity-wise, not wise!

However: the dual formulation is entirely expressible in terms of the kernel function  $k(\cdot, \cdot)$

6

If we use the kernel trick (if it is possible), we can work directly with kernels, and avoid the explicit introduction of the feature transformation  $\phi(\underline{x})$ . This allows us to use features of high (even infinite) dimensionality.

## CONSTRUCTING KERNEL FUNCTIONS DIRECTLY

Example:  $k(\underline{x}, \underline{z}) \triangleq (\underline{x}^T \underline{z})^2$

Recap:  $k(\underline{x}, \underline{x}') = \underbrace{\phi^T(\underline{x}) \phi(\underline{x}')}_{\substack{\text{kernel} \\ \text{(scalar)}}}$

$\phi$ : mapping function

original space

$$\underline{x}^T \underline{z} = [x_2 \ x_1] \begin{bmatrix} z_2 \\ z_1 \end{bmatrix} = x_2 z_2 + x_1 z_1$$

$$\Rightarrow (\underline{x}^T \underline{z})^2 = (x_2 z_2 + x_1 z_1)^2 = x_2^2 z_2^2 + 2 x_1 z_1 x_2 z_2 + x_1^2 z_1^2$$

try to separate into

$$\begin{array}{c} \boxed{\phi(\underline{x})} \\ \downarrow \\ 2\text{-D} \end{array} \quad \begin{array}{c} \boxed{\phi(\underline{z})} \\ \downarrow \\ 2\text{-D} \end{array}$$

7

$$\begin{bmatrix} x_2^2 & \sqrt{2} x_2 x_1 & x_1^2 \end{bmatrix} \begin{bmatrix} x_2^2 \\ \sqrt{2} x_2 x_1 \\ x_1^2 \end{bmatrix}$$

$$\phi(\underline{x}) = \begin{bmatrix} x_2^2 \\ \sqrt{2} x_2 x_1 \\ x_1^2 \end{bmatrix}$$

$\downarrow$   
 $\begin{bmatrix} x_2 \\ x_1 \end{bmatrix}$   
 (2-D)

(3-D)

mapping, comprises all second order terms with a specific weighing between them.