

(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  $\uparrow$  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

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

$N \times 1$

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

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

[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

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

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

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

typically,  $N \gg M$

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

$N \times N$  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:  $k(\underline{x}, \underline{x}') = \underbrace{\phi^T(\underline{x})}_{\text{kernel (scalar)}} \underbrace{\phi(\underline{x}')}_{\phi: \text{mapping function}}$

2-D

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

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

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.

(Mercer's Condition) Necessary & Sufficient condition for a function  $k(\underline{x}, \underline{x}')$  to be a valid kernel:

The Gram Matrix  $K$   $K(i, j) = k(\underline{x}_i, \underline{x}_j)$  should be PSD  $\forall \underline{x}_i$ :

$$\text{i.e., } \underline{x}_i^T K \underline{x}_i \geq 0$$

Properties: given valid kernels  $k_1(\underline{x}, \underline{x}')$  and  $k_2(\underline{x}, \underline{x}')$ , the following new kernels will also be valid kernels:-

$$1) k(\underline{x}, \underline{x}') = c k_1(\underline{x}, \underline{x}'), \quad c > 0$$

$$2) k(\underline{x}, \underline{x}') = f(\underline{x}) k_1(\underline{x}, \underline{x}') f(\underline{x}')$$

$$3) k(\underline{x}, \underline{x}') = q(k_1(\underline{x}, \underline{x}')), \quad \text{where } q(\cdot) \text{ is a function polynomial with non-negative coefficients}$$

$$4) k(\underline{x}, \underline{x}') = \exp(k_1(\underline{x}, \underline{x}'))$$

$$5) k(\underline{x}, \underline{x}') = k_1(\underline{x}, \underline{x}') + k_2(\underline{x}, \underline{x}')$$

$$6) k(\underline{x}, \underline{x}') = k_1(\underline{x}, \underline{x}') k_2(\underline{x}, \underline{x}')$$

$$7) k(\underline{x}, \underline{x}') = k_3(\phi(\underline{x}), \phi(\underline{x}')), \quad \phi(\underline{x}): \underline{x} \rightarrow \mathbb{R}^M$$

$$8) k(\underline{x}, \underline{x}') = k_a(\underline{x}_a, \underline{x}'_a) + k_b(\underline{x}_b, \underline{x}'_b)$$

$\underline{x} = (\underline{x}_a, \underline{x}_b)$ :  $\underline{x}_a$  &  $\underline{x}_b$  are not necessarily disjoint,  $k_a(\cdot)$  and  $k_b(\cdot)$  are valid kernel functions.

$$9) k(\underline{x}, \underline{x}') = k_a(\underline{x}_a, \underline{x}'_a) k_b(\underline{x}_b, \underline{x}'_b)$$

# "Gaussian kernel"

Interesting Example:  $k(\underline{x}, \underline{x}') = \exp\left(-\frac{\|\underline{x} - \underline{x}'\|^2}{2\sigma^2}\right)$

$$\begin{aligned}\|\underline{x} - \underline{x}'\|^2 &= (\underline{x} - \underline{x}')^T (\underline{x} - \underline{x}') \\ &= (\underline{x}^T - \underline{x}'^T) (\underline{x} - \underline{x}') \\ &= \underline{x}^T \underline{x} + \underline{x}'^T \underline{x}' - 2\underline{x}^T \underline{x}'\end{aligned}$$

$$k(\underline{x}, \underline{x}') = \underbrace{\exp\left(-\frac{\underline{x}^T \underline{x}}{2\sigma^2}\right)}_{\text{exp}\left(\frac{-\underline{x}^T \underline{x}}{2\sigma^2}\right)} \cdot \underbrace{\exp\left(-\frac{\underline{x}'^T \underline{x}'}{2\sigma^2}\right)}_{\text{exp}\left(\frac{-\underline{x}'^T \underline{x}'}{2\sigma^2}\right)} \cdot \underbrace{\exp\left(\frac{\underline{x}^T \underline{x}'}{\sigma^2}\right)}_{\text{exp}\left(\frac{\underline{x}^T \underline{x}'}{\sigma^2}\right)}$$

$$\underbrace{\exp\left(\frac{-\underline{x}^T \underline{x}}{2\sigma^2}\right) \exp\left(\frac{-\underline{x}'^T \underline{x}'}{2\sigma^2}\right)}_{\text{exp}\left(\frac{-\underline{x}^T \underline{x} - \underline{x}'^T \underline{x}'}{2\sigma^2}\right)} \exp\left(\frac{\underline{x}^T \underline{x}'}{\sigma^2}\right)$$

$\underline{x}^T \underline{x}'$  is a kernel (linear kernel)

$\Rightarrow \left(\frac{1}{\sigma} \underline{x}^T\right) \left(\frac{1}{\sigma} \underline{x}'\right)$  is also a kernel

$\Rightarrow \frac{\underline{x}^T \underline{x}'}{\sigma^2}$  is a kernel

$\Rightarrow \exp\left(\frac{\underline{x}^T \underline{x}'}{\sigma^2}\right)$  is also a kernel

$\underbrace{\exp\left(\frac{\underline{x}^T \underline{x}'}{\sigma^2}\right)}_{k_1(\underline{x}, \underline{x}')} \quad \left(\because \exp(k_1(\underline{x}, \underline{x}')) \text{ is a kernel}\right)$

$\Rightarrow f(\underline{x}) k_1(\underline{x}, \underline{x}') \cdot f(\underline{x}')$  is also a kernel

$$\exp\left(\frac{-\underline{x}^T \underline{x}}{2\sigma^2}\right)$$

$$\exp\left(\frac{-\underline{x}'^T \underline{x}'}{2\sigma^2}\right)$$

Q.E.D.