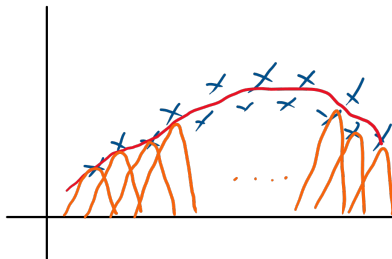
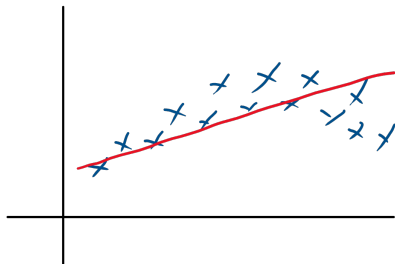


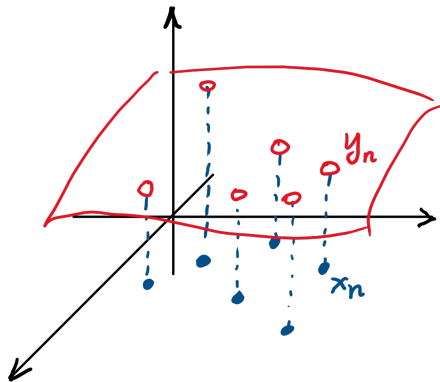
## Regression with Kernels \*

## Why Another Method?

- Linear regression: Pick a **global** model, best fit globally.
- Kernel method: Pick a **local** model, best fit locally.
- In kernel method, instead of picking a line / a quadratic equation, we pick a **kernel**.
- A kernel is a measure of **distance** between **training samples**.
- Kernel method buys us the ability to handle nonlinearity.
- Ordinary regression is based on the **columns** (features) of  $\mathbf{A}$ .
- Kernel method is based on the **rows** (samples) of  $\mathbf{A}$ .



# Pictorial Illustration



goal: learn the surface

prediction: When new  
sample comes, interpolate  
on the surface

# Overview of the Method

## Model Parameter:

- We **want** the model parameter  $\hat{\theta}$  to look like: (How? Question 1)

$$\hat{\theta} = \sum_{n=1}^N \alpha_n \mathbf{x}^n.$$

- This model expresses  $\hat{\theta}$  as a combination of the **samples**.
- The trainable parameters are  $\alpha_n$ , where  $n = 1, \dots, N$ .
- If we can make  $\alpha_n$  **local**, i.e., non-zero for only a few of them, then we can achieve our goal: localized, sample-dependent.

## Predicted Value

- The predicted value of a new sample  $\mathbf{x}$  is

$$\hat{y} = \hat{\theta}^T \mathbf{x} = \sum_{n=1}^N \alpha_n \langle \mathbf{x}, \mathbf{x}^n \rangle.$$

- We **want** this model to encapsulate nonlinearity. (How? Question 2)

## Dual Form of Linear Regression

**Goal:** Addresses Question 1: Express  $\hat{\boldsymbol{\theta}}$  as

$$\hat{\boldsymbol{\theta}} = \sum_{n=1}^N \alpha_n \mathbf{x}^n.$$

We start by listing out a technical lemma:

### Lemma

For any  $\mathbf{A} \in \mathbb{R}^{N \times d}$ ,  $\mathbf{y} \in \mathbb{R}^d$ , and  $\lambda > 0$ ,

$$(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{y}. \quad (1)$$

Proof: See Appendix.

Remark:

- The dimensions of  $\mathbf{I}$  on the left is  $d \times d$ , on the right is  $N \times N$ .
- If  $\lambda = 0$ , then the above is true only when  $\mathbf{A}$  is invertible.

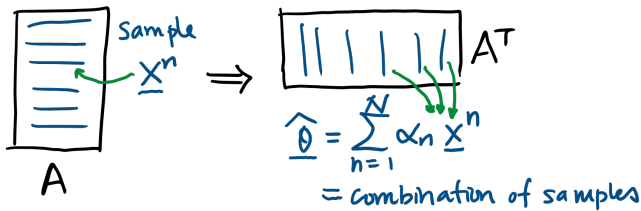
# Dual Form of Linear Regression

- Using the Lemma, we can show that

$$\hat{\theta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y} \quad (\text{Primal Form})$$

$$= \mathbf{A}^T \underbrace{(\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{y}}_{\stackrel{\text{def}}{=} \boldsymbol{\alpha}} \quad (\text{Dual Form})$$

$$= \begin{bmatrix} - & (\mathbf{x}^1)^T & - \\ - & (\mathbf{x}^2)^T & - \\ & \vdots & \\ - & (\mathbf{x}^N)^T & - \end{bmatrix}^T \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_N \end{bmatrix} = \sum_{n=1}^N \alpha_n \mathbf{x}^n, \quad \alpha_n \stackrel{\text{def}}{=} [(\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{y}]_n$$



# The Kernel Trick

**Goal:** Addresses Question 2: Introduce nonlinearity to

$$\hat{y} = \hat{\boldsymbol{\theta}}^T \mathbf{x} = \sum_{n=1}^N \alpha_n \langle \mathbf{x}, \mathbf{x}^n \rangle.$$

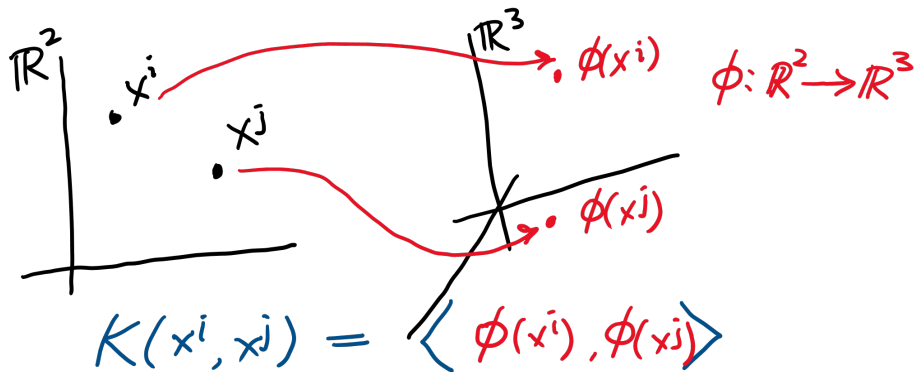
**The Idea:**

- Replace the inner product  $\langle \mathbf{x}, \mathbf{x}^n \rangle$  by  $k(\mathbf{x}, \mathbf{x}^n)$ :

$$\hat{y} = \hat{\boldsymbol{\theta}}^T \mathbf{x} = \sum_{n=1}^N \alpha_n k(\mathbf{x}, \mathbf{x}^n).$$

- $k(\cdot, \cdot)$  is called a **kernel**.
- A kernel is a measure of the **distance** between two samples  $\mathbf{x}^i$  and  $\mathbf{x}^j$ .
- $\langle \mathbf{x}^i, \mathbf{x}^j \rangle$  measure distance in the ambient space,  $k(\mathbf{x}^i, \mathbf{x}^j)$  measure distance in a **transformed** space.
- In particular, a valid kernel takes the form  $k(\mathbf{x}^i, \mathbf{x}^j) = \langle \phi(\mathbf{x}^i), \phi(\mathbf{x}^j) \rangle$  for some nonlinear transforms  $\phi$ .

## Kernels Illustrated



- A kernel typically lifts the ambient dimension to a **higher** one.
- For example, mapping from  $\mathbb{R}^2$  to  $\mathbb{R}^3$

$$\mathbf{x}^n = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \text{and} \quad \phi(\mathbf{x}_n) = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{bmatrix}$$



## Relationship between Kernel and Transform

Consider the following kernel  $k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v})^2$ . What is the transform?

- Suppose  $\mathbf{u}$  and  $\mathbf{v}$  are in  $\mathbb{R}^2$ . Then  $(\mathbf{u}^T \mathbf{v})^2$  is

$$\begin{aligned}(\mathbf{u}^T \mathbf{v})^2 &= \left( \sum_{i=1}^2 u_i v_i \right) \left( \sum_{j=1}^2 u_j v_j \right) \\&= \sum_{i=1}^2 \sum_{j=1}^2 (u_i u_j)(v_i v_j) = \begin{bmatrix} u_1^2 & u_1 u_2 & u_2 u_1 & u_2^2 \end{bmatrix} \begin{bmatrix} v_1^2 \\ v_1 v_2 \\ v_2 v_1 \\ v_2^2 \end{bmatrix}.\end{aligned}$$

- So if we define  $\phi$  as

$$\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \mapsto \phi(\mathbf{u}) = \begin{bmatrix} u_1^2 \\ u_1 u_2 \\ u_2 u_1 \\ u_2^2 \end{bmatrix}$$

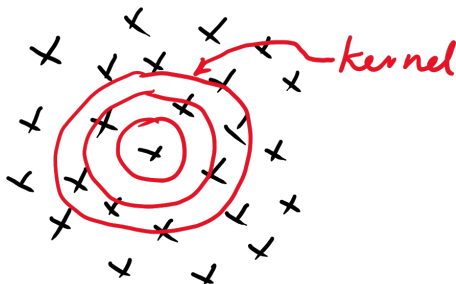
then  $(\mathbf{u}^T \mathbf{v})^2 = \langle \phi(\mathbf{u}), \phi(\mathbf{v}) \rangle$ .

# Radial Basis Function

A useful kernel is the **radial basis kernel** (RBF):

$$k(\mathbf{u}, \mathbf{v}) = \exp \left\{ -\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2} \right\}.$$

- The corresponding nonlinear transform of RBF is **infinite dimensional**. See Appendix.
- $\|\mathbf{u} - \mathbf{v}\|^2$  measures the distance between two data points  $\mathbf{u}$  and  $\mathbf{v}$ .
- $\sigma$  is the std dev, defining “far” and “close”.
- RBF enforces **local** structure; Only a few samples are used.



# Kernel Method

Given the choice of the kernel function, we can write down the algorithm as follows.

- 1 Pick a kernel function  $k(\cdot, \cdot)$ .
- 2 Construct a kernel matrix  $\mathbf{K} \in \mathbb{R}^{N \times N}$ , where  $[\mathbf{K}]_{ij} = k(\mathbf{x}^i, \mathbf{x}^j)$ , for  $i = 1, \dots, N$  and  $j = 1, \dots, N$ .
- 3 Compute the coefficients  $\boldsymbol{\alpha} \in \mathbb{R}^N$ , with

$$\alpha_n = [(\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}]_n.$$

- 4 Estimate the predicted value for a new sample  $\mathbf{x}$ :

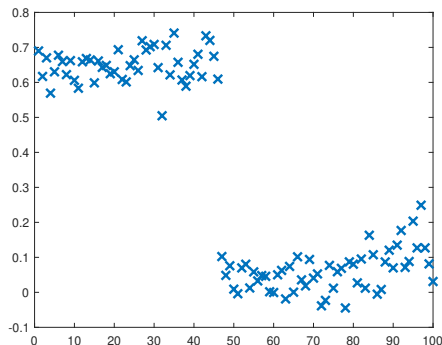
$$g_{\theta}(\mathbf{x}) = \sum_{n=1}^N \alpha_n k(\mathbf{x}, \mathbf{x}^n).$$

Therefore, the choice of the regression function is shifted to the choice of the kernel.

## Example

**Goal:** Use the kernel method to fit the data points shown as follows.

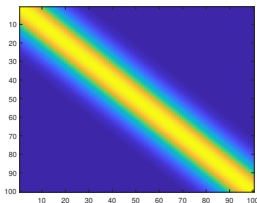
- What is the input feature vector  $\mathbf{x}^n$ ?  $\mathbf{x}^n = t_n$ : The time stamps.
- What is the output  $y_n$ ?  $y^n$  is the height.
- Which kernel to choose? Let us consider the RBF.



## Example (using RBF)

- Define the fitted function as  $g_{\theta}(t)$ . [Here,  $\theta$  refers to  $\alpha$ .]
- The RBF is defined as  $k(t_i, t_j) = \exp\{-(t_i - t_j)^2/2\sigma^2\}$ , for some  $\sigma$ .
- The matrix  $\mathbf{K}$  looks something below

$$[\mathbf{K}]_{ij} = \exp\{-(t_i - t_j)^2/2\sigma^2\}.$$



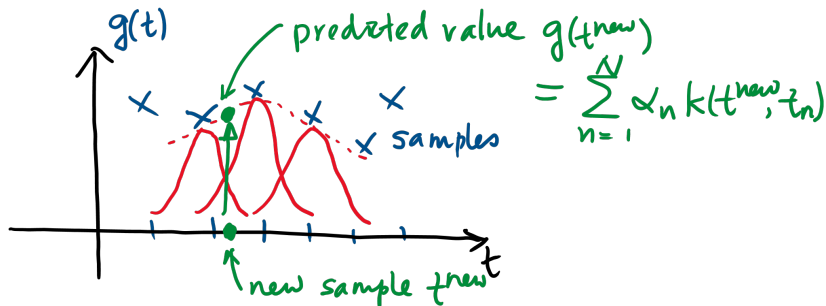
- $\mathbf{K}$  is a banded diagonal matrix. (Why?)
- The coefficient vector is  $\alpha_n = [(\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}]_n$ .

## Example (using RBF)

- Using the RBF, the predicted value is given by

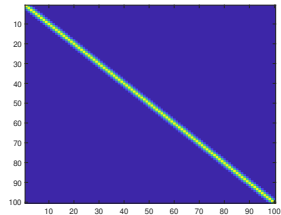
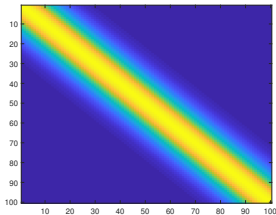
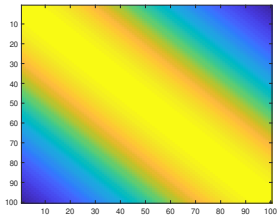
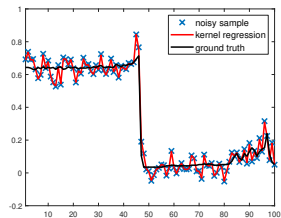
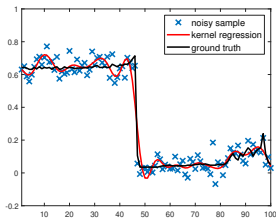
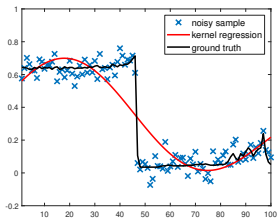
$$g_{\theta}(t^{\text{new}}) = \sum_{n=1}^N \alpha_n k(t^{\text{new}}, t_n) = \sum_{n=1}^N \alpha_n e^{-\frac{(t^{\text{new}} - t_n)^2}{2\sigma^2}}.$$

- Pictorially, the predicted function  $g_{\theta}$  can be viewed as the linear combination of the Gaussian kernels.



## Effect of $\sigma$

- Large  $\sigma$ : Flat kernel. Over-smoothing.
- Small  $\sigma$ : Narrow kernel. Under-smoothing.
- Below shows an example of the fitting and the kernel matrix  $K$ .



Too large

About right

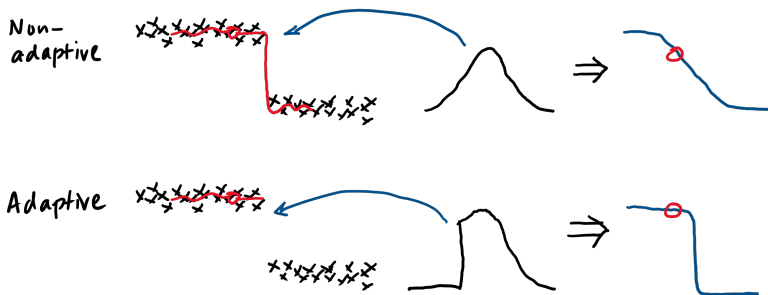
Too small

# Any Improvement?

- We can improve the above kernel by considering  $\mathbf{x}^n = [y_n, t_n]^T$ .
- Define the kernel as

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp \left\{ - \left( \frac{(y_i - y_j)^2}{2\sigma_r^2} + \frac{(t_i - t_j)^2}{2\sigma_s^2} \right) \right\}.$$

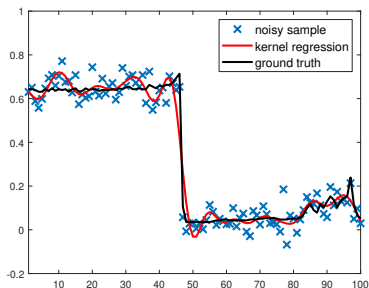
- This new kernel is adaptive (**edge-aware**).



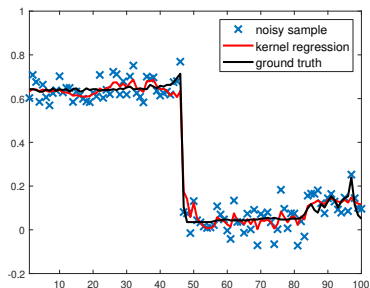


# Any Improvement?

Here is a comparison.



without improvement

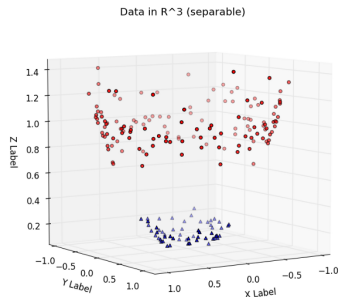
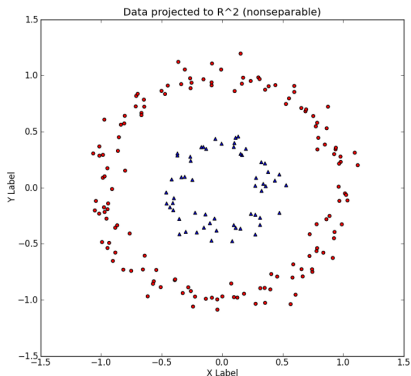


with improvement

- This idea is known as **bilateral filter** in the computer vision literature.
- Can be further extended to 2D image where  $\mathbf{x}^n = [y_n, \mathbf{s}_n]$ , for some spatial coordinate  $\mathbf{s}_n$ .
-

# Kernel Methods in Classification

- The concept of lifting the data to higher dimension is useful for classification.<sup>1</sup>

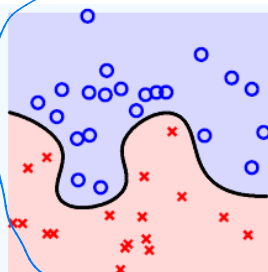


<sup>1</sup>Image source:

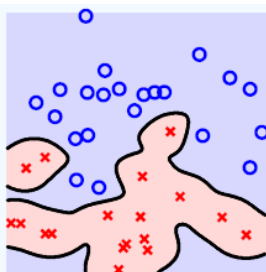
# Kernels in Support Vector Machines

**Example.** RBF for SVM (We will discuss SVM later in the semester.)

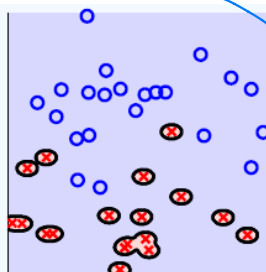
- Radial Basis Function is often used in support vector machine.
- Poor choice of parameter can lead to low training loss, but with the risk of over-fit.
- Under-fitted data can sometimes give better generalization.



$$\exp(-1\|x - x'\|^2)$$



$$\exp(-10\|x - x'\|^2)$$



$$\exp(-100\|x - x'\|^2)$$

## **Appendix**

# Proof of Lemma

## Lemma

For any matrix  $\mathbf{A} \in \mathbb{R}^{N \times d}$ ,  $\mathbf{y} \in \mathbb{R}^d$ , and  $\lambda > 0$ ,

$$(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{y}. \quad (2)$$

- The left hand side is solution to normal equation, which means  $\mathbf{A}^T \mathbf{A} \boldsymbol{\theta} + \lambda \boldsymbol{\theta} = \mathbf{A}^T \mathbf{y}$ .
- Rearrange terms gives  $\boldsymbol{\theta} = \mathbf{A}^T \left[ \frac{1}{\lambda} (\mathbf{y} - \mathbf{A} \boldsymbol{\theta}) \right]$ .
- Define  $\boldsymbol{\alpha} = \frac{1}{\lambda} (\mathbf{y} - \mathbf{A} \boldsymbol{\theta})$ , then  $\boldsymbol{\theta} = \mathbf{A}^T \boldsymbol{\alpha}$ .
- Substitute  $\boldsymbol{\theta} = \mathbf{A}^T \boldsymbol{\alpha}$  into  $\boldsymbol{\alpha} = \frac{1}{\lambda} (\mathbf{y} - \mathbf{A} \boldsymbol{\theta})$ , we have

$$\boldsymbol{\alpha} = \frac{1}{\lambda} (\mathbf{y} - \mathbf{A} \mathbf{A}^T \boldsymbol{\alpha}).$$

- Rearrange terms gives  $(\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I}) \boldsymbol{\alpha} = \mathbf{y}$ , which yields  $\boldsymbol{\alpha} = (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$ .
- Substitute into  $\boldsymbol{\theta} = \mathbf{A}^T \boldsymbol{\alpha}$  gives  $\boldsymbol{\theta} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$ .

## Non-Linear Transform for RBF

- Let us consider scalar  $u \in \mathbb{R}$ .

$$\begin{aligned}k(u, v) &= \exp\{-(u - v)^2\} \\&= \exp\{-u^2\} \exp\{2uv\} \exp\{-v^2\} \\&= \exp\{-u^2\} \left( \sum_{k=0}^{\infty} \frac{2^k u^k v^k}{k!} \right) \exp\{-v^2\} \\&= \exp\{-u^2\} \left( 1, \sqrt{\frac{2^1}{1!}} u, \sqrt{\frac{2^2}{2!}} u^2, \sqrt{\frac{2^3}{3!}} u^3, \dots, \right)^T \\&\quad \times \left( 1, \sqrt{\frac{2^1}{1!}} v, \sqrt{\frac{2^2}{2!}} v^2, \sqrt{\frac{2^3}{3!}} v^3, \dots, \right) \exp\{-v^2\}\end{aligned}$$

- So  $\Phi$  is

$$\phi(x) = \exp\{-x^2\} \left( 1, \sqrt{\frac{2^1}{1!}} x, \sqrt{\frac{2^2}{2!}} x^2, \sqrt{\frac{2^3}{3!}} x^3, \dots, \right)$$

## Kernels are Positive Semi-Definite

Given  $\{\mathbf{x}_j\}_{j=1}^N$ , construct a  $N \times N$  matrix  $\mathbf{K}$  such that

$$[\mathbf{K}]_{ij} = K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j).$$

**Claim:**  $\mathbf{K}$  is positive semi-definite.

Let  $\mathbf{z}$  be an arbitrary vector. Then,

$$\begin{aligned} \mathbf{z}^T \mathbf{K} \mathbf{z} &= \sum_{i=1}^n \sum_{j=1}^N z_i K_{ij} z_j = \sum_{i=1}^N \sum_{j=1}^N z_i \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j) z_j \\ &= \sum_{i=1}^N \sum_{j=1}^N z_i \left( \sum_{k=1}^N [\Phi(\mathbf{x}_i)]_k [\Phi(\mathbf{x}_j)]_k \right) z_j \stackrel{(a)}{=} \sum_{k=1}^N \left( \sum_{i=1}^N [\Phi(\mathbf{x}_i)]_k z_i \right)^2 \geq 0 \end{aligned}$$

where  $[\Phi(\mathbf{x}_i)]_k$  denotes the  $k$ -th element of the vector  $\Phi(\mathbf{x}_i)$ .

# Existence of Nonlinear Transform

- We just showed that: If  $K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)$  for any  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , then  $\mathbf{K}$  is symmetric positive semi-definite.
- The converse also holds: If  $\mathbf{K}$  is symmetric positive semi-definite for any  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , then there exist  $\Phi$  such that  $K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)$ .
- This converse is difficult to prove.
- It is called the **Mercer Condition**.
- Kernels satisfying Mercer's condition have  $\Phi$ .
- You can use the condition to rule out invalid kernels.
- But proving a valid kernel is still hard.