

# ECE595 / STAT598: Machine Learning I

## Lecture 3.1: Regression with Kernels - Kernel Method

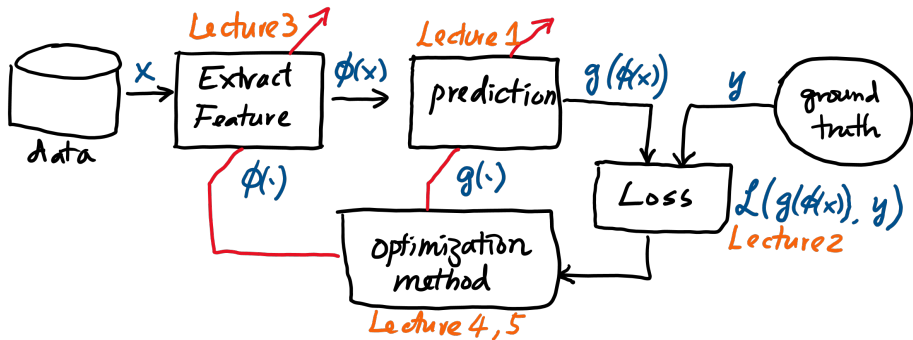
Spring 2020

Stanley Chan

School of Electrical and Computer Engineering  
Purdue University



# Outline



# Outline

## Mathematical Background

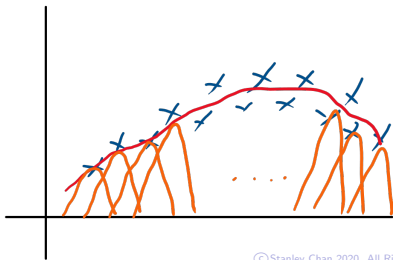
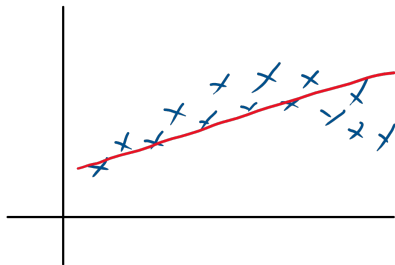
- Lecture 1: Linear regression: A basic data analytic tool
- Lecture 2: Regularization: Constraining the solution
- **Lecture 3: Kernel Method: Enabling nonlinearity**

## Lecture 3: Kernel Method

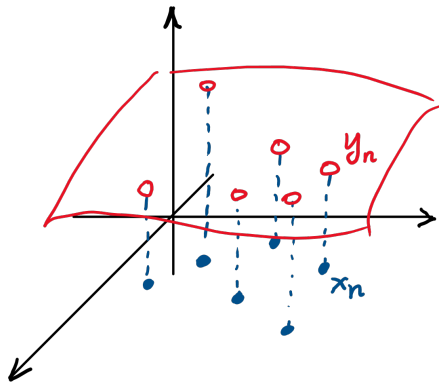
- **Kernel Method**
  - **Dual Form**
  - **Kernel Trick**
  - **Algorithm**
- **Examples**
  - Radial Basis Function (RBF)
  - Regression using RBF
  - Kernel Methods in Classification

## 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.

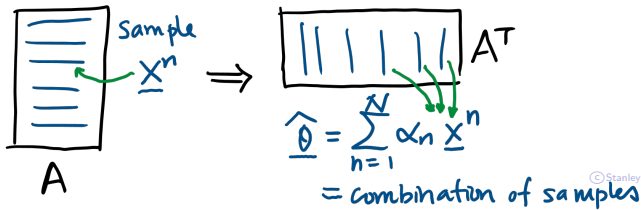
# Dual Form of Linear Regression

- Using the Lemma, we can show that

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

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

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





# The Kernel Trick

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

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

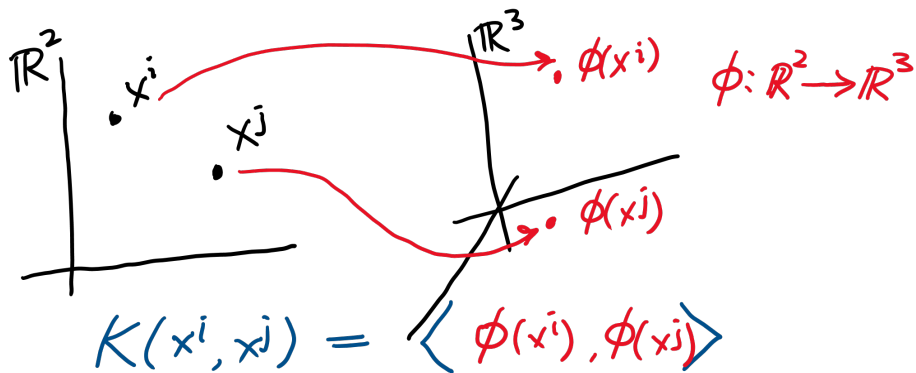
**The Idea:**

- 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

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

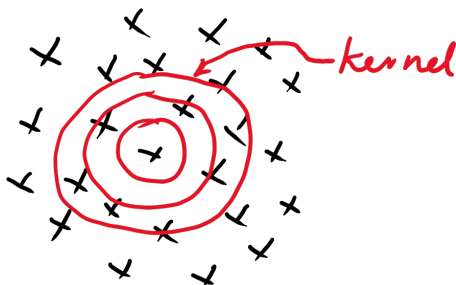
then  $(\mathbf{u}^T \mathbf{v})^2 = \langle \phi(\mathbf{u}), \phi(\mathbf{v}) \rangle$ .

# 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  $\mathbf{x}$ :

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

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