# ECE595 / STAT598: Machine Learning I Lecture 03: Regression with Kernels

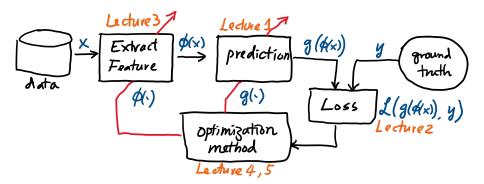
Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



## Outline



### Outline

#### Mathematical Background

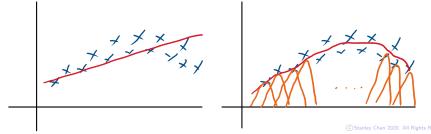
- Lecture 1: Linear regression: A basic data analytic tool
- Lecture 2: Regularization: Constraining the solution
- Lecture 3: Kernel Method: Enabling nonlinearity

#### Lecture 3: Kernel Method

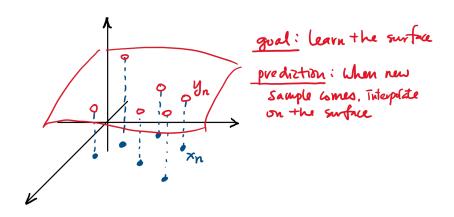
- Kernel Method
  - Dual Form
  - Kernel Trick
  - Algorithm
- Examples
  - Radial Basis Function (RBF)
  - Regression using RBF
  - Kernel Methods in Classification

# 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 A.
- Kernel method is based on the rows (samples) of A.



## Pictorial Illustration



## Overview of the Method

#### **Model Parameter:**

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

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

- This model expresses  $\widehat{\theta}$  as a combination of the samples.
- The trainable parameters are  $\alpha_n$ , where n = 1, ..., 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 x is

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

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

# Dual Form of Linear Regression

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

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

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  $\lambda > 0$ ,

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

Proof: See Appendix.

#### Remark:

- The dimensions of 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  ${\it A}$  is invertible.

# Dual Form of Linear Regression

• Using the Lemma, we can show that 
$$\widehat{\theta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y} \qquad \text{(Primal Form)}$$

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

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

### The Kernel Trick

Goal: Addresses Question 2: Introduce nonlinearity to

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

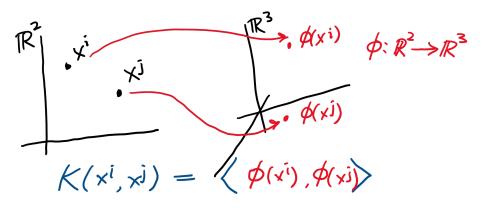
#### The Idea:

• Replace the inner product  $\langle x, x^n \rangle$  by  $k(x, x^n)$ :

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

- $k(\cdot, \cdot)$  is called a **kernel**.
- A kernel is a measure of the **distance** between two samples  $x^i$  and  $x^j$ .
- $\langle x^i, x^j \rangle$  measure distance in the ambient space,  $k(x^i, x^j)$  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  $\mathbb{R}^3$

$$m{x}^n = egin{bmatrix} x_1 \ x_2 \end{bmatrix}$$
 and  $\phi(m{x}_n) = egin{bmatrix} x_1^2 \ x_1 x_2 \ x_2^2 \end{bmatrix}$  © Stanley Chan 2020. All Rights Reserves 10 / 28

# 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  $\boldsymbol{u}$  and  $\boldsymbol{v}$  are in  $\mathbb{R}^2$ . Then  $(\boldsymbol{u}^T\boldsymbol{v})^2$  is

$$(\mathbf{u}^{T}\mathbf{v})^{2} = \left(\sum_{i=1}^{2} u_{i}v_{i}\right) \left(\sum_{j=1}^{2} u_{j}v_{j}\right)$$

$$= \sum_{i=1}^{2} \sum_{j=1}^{2} (u_{i}u_{j})(v_{i}v_{j}) = \begin{bmatrix} u_{1}^{2} & u_{1}u_{2} & u_{2}u_{1} & u_{2}^{2} \end{bmatrix} \begin{bmatrix} v_{1}^{2} \\ v_{1}v_{2} \\ v_{2}v_{1} \\ v_{2}^{2} \end{bmatrix}.$$

• So if we define  $\phi$  as

$$m{u} = egin{bmatrix} u_1 \ u_2 \end{bmatrix} \quad \mapsto \quad \phi(m{u}) = egin{bmatrix} u_1^2 \ u_1 u_2 \ u_2 u_1 \ u_2^2 \end{bmatrix}$$

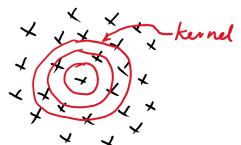
then  $(\mathbf{u}^T \mathbf{v})^2 = \langle \phi(\mathbf{u}), \phi(\mathbf{v}) \rangle$ .

#### 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  $\mathbf{v}$ .
- $\bullet$   $\sigma$  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, ..., N and j = 1, ..., N.
- **3** Compute the coefficients  $\alpha \in \mathbb{R}^N$ , with

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

**\bullet** Estimate the predicted value for a new sample 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.

### Outline

#### Mathematical Background

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

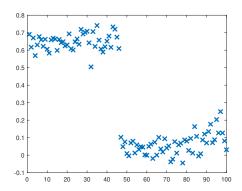
#### Lecture 3: Kernel Method

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

## Example

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

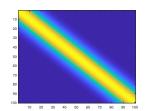
- What is the input feature vector  $\mathbf{x}^n$ ?  $\mathbf{x}^n = t_n$ : The time stamps.
- What is the output  $y_n$ ?  $y^n$  is the height.
- Which kernel to choose? Let us consider the RBF.



# Example (using RBF)

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

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



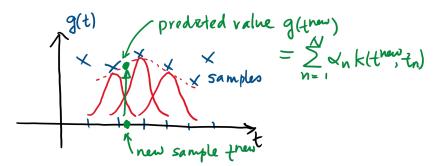
- K is a banded diagonal matrix. (Why?)
- The coefficient vector is  $\alpha_n = [(\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}]_n$ .

# Example (using RBF)

Using the RBF, the predicted value is given by

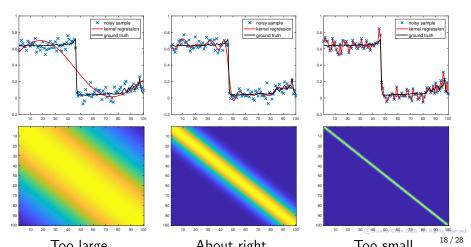
$$g_{\theta}(t^{\text{new}}) = \sum_{n=1}^{N} \alpha_n k(t^{\text{new}}, t_n) = \sum_{n=1}^{N} \alpha_n e^{-\frac{(t^{\text{new}} - t_n)^2}{2\sigma^2}}.$$

• 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.
- ullet Below shows an example of the fitting and the kernel matrix  $oldsymbol{\mathit{K}}$ .

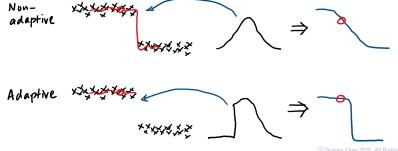


# Any Improvement?

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

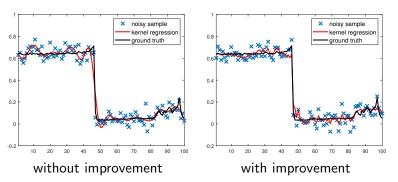
$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp \left\{ -\left( \frac{(y_i - y_j)^2}{2\sigma_r^2} + \frac{(t_i - t_j)^2}{2\sigma_s^2} \right) \right\}.$$

This new kernel is adaptive (edge-aware).



# Any Improvement?

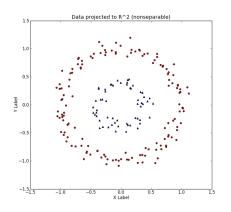
Here is a comparison.

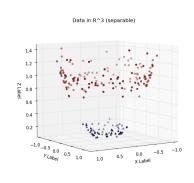


- This idea is known as **bilateral filter** in the computer vision literature.
- Can be further extended to 2D image where  $\mathbf{x}^n = [y_n, \mathbf{s}_n]$ , for some spatial coordinate  $\mathbf{s}_n$ .
- Many applications. See Reading List.

### Kernel Methods in Classification

 The concept of lifting the data to higher dimension is useful for classification.



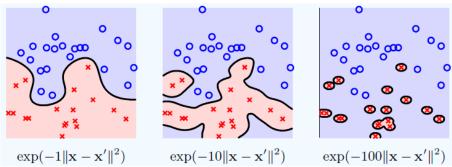


<sup>&</sup>lt;sup>1</sup>Image source:

# 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.



# Reading List

#### Kernel Method:

- Learning from Data (Chapter 3.4)
   https://work.caltech.edu/telecourse
- CMU 10701 Lecture 4 https://www.cs.cmu.edu/~tom/10701\_sp11/slides/Kernels\_SVM\_04\_7\_2011-ann.pdf
- Berkeley CS 194 Lecture 7 https: //people.eecs.berkeley.edu/~russell/classes/cs194/
- Oxford C19 Lecture 3
   http://www.robots.ox.ac.uk/~az/lectures/ml/lect3.pdf

## Kernel Regression in Computer Vision:

- Bilateral Filter https://people.csail.mit.edu/sparis/bf\_ course/course\_notes.pdf
- Takeda and Milanfar, "Kernel regression for image processing and reconstruction", IEEE Trans. Image Process. (2007) https://ieeexplore.ieee.org/document/4060955

**Appendix** 

#### Proof of Lemma

#### Lemma

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

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

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

$$oldsymbol{lpha} = rac{1}{\lambda} (oldsymbol{y} - oldsymbol{A} oldsymbol{A}^T oldsymbol{lpha}).$$

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

## Non-Linear Transform for RBF

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

$$k(u, v) = \exp\{-(u - v)^2\}$$

$$= \exp\{-u^2\} \exp\{2uv\} \exp\{-v^2\}$$

$$= \exp\{-u^2\} \left(\sum_{k=0}^{\infty} \frac{2^k u^k v^k}{k!}\right) \exp\{-v^2\}$$

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

$$\times \left(1, \sqrt{\frac{2^1}{1!}} v, \sqrt{\frac{2^2}{2!}} v^2, \sqrt{\frac{2^3}{3!}} v^3, \dots, \right) \exp\{-v^2\}$$

So Φ 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  $\{x_j\}_{j=1}^N$ , construct a  $N \times N$  matrix K such that

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

Claim: K is positive semi-definite.

Let z be an arbitrary vector. Then,

$$z^{T}Kz = \sum_{i=1}^{n} \sum_{j=1}^{N} z_{i}K_{ij}z_{j} = \sum_{i=1}^{N} \sum_{j=1}^{N} z_{i}\Phi(x_{i})^{T}\Phi(x_{j})z_{j}$$

$$= \sum_{i=1}^{N} \sum_{j=1}^{N} z_{i} \left(\sum_{k=1}^{N} [\Phi(x_{i})]_{k} [\Phi(x_{j})]_{k}\right) z_{j} \stackrel{\text{(a)}}{=} \sum_{k=1}^{N} \left(\sum_{j=1}^{N} [\Phi(x_{j})]_{k}z_{j}\right)^{2} \ge 0$$

where  $[\Phi(x_i)]_k$  denotes the k-th element of the vector  $\Phi(x_i)$ .

## Existence of Nonlinear Transform

- We just showed that: If  $K(x_i, x_j) = \Phi(x_i)^T \Phi(x_j)$  for any  $x_1, \dots, x_N$ , then K is symmetric positive semi-definite.
- The converse also holds: If K is symmetric positive semi-definite for any  $x_1, \ldots, x_N$ , then there exist  $\Phi$  such that  $K(x_i, x_i) = \Phi(x_i)^T \Phi(x_i)$ .
- This converse is difficult to prove.
- It is called the Mercer Condition.
- Kernels satisfying Mercer's condition have Φ.
- You can use the condition to rule out invalid kernels.
- But proving a valid kernel is still hard.