

# Pattern Information Processing: Linear Models and Least-Squares<sup>1</sup>

Masashi Sugiyama  
(Department of Computer Science)

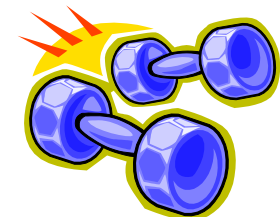
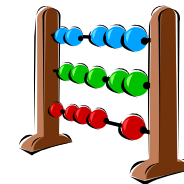
Contact: W8E-505

[sugi@cs.titech.ac.jp](mailto:sugi@cs.titech.ac.jp)

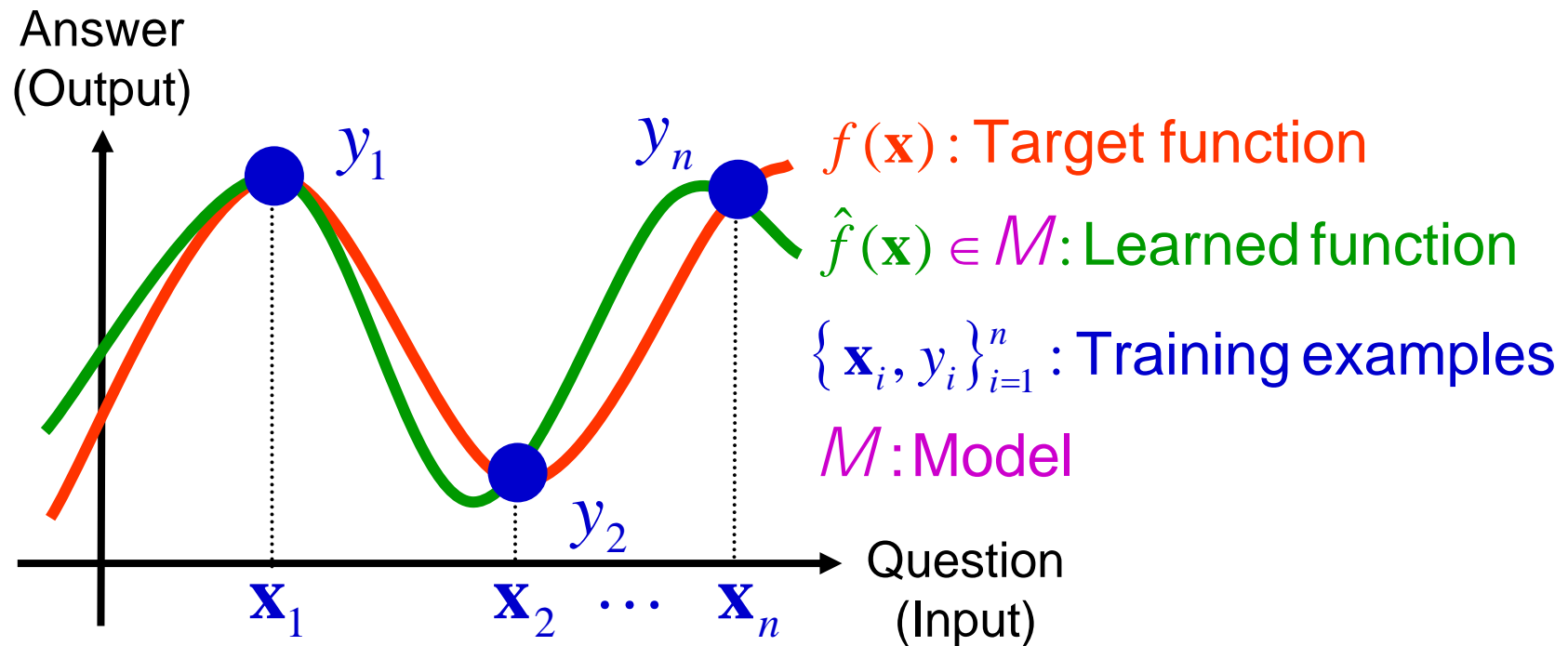
<http://sugiyama-www.cs.titech.ac.jp/~sugi/>

# Focus of This Course

- There are 3 topics in learning research.
  - Understanding human brains
  - Developing learning machines
  - Mathematically clarifying mechanism of learning
- There are 3 types of learning.
  - Supervised learning
  - Unsupervised learning
  - Reinforcement learning
- Topics of supervised learning:
  - Active learning
  - Model selection
  - Learning methods



# Supervised Learning As Function Approximation



Using training examples  $\{\mathbf{x}_i, y_i\}_{i=1}^n$ ,  
find a function  $\hat{f}(\mathbf{x})$  from a model  $\mathcal{M}$   
that well approximates the target function  $f(\mathbf{x})$ .

# Formal Notation and Assumptions<sup>4</sup>

- $f(\mathbf{x})$  : Learning target function
- $\mathcal{D} \subset \mathbb{R}^d$  : Domain of  $f(\mathbf{x})$
- $\mathbf{x}_i$  : Training input point
- $y_i = f(\mathbf{x}_i) + \epsilon_i$  : Training output value
- $\epsilon_i$  : mean zero, independent and identically distributed (“i.i.d.”)  
$$\mathbb{E}_{\epsilon}[\epsilon_i] = 0 \quad \mathbb{E}_{\epsilon}[\epsilon_i \epsilon_j] = \begin{cases} \sigma^2 & (i = j) \\ 0 & (i \neq j) \end{cases}$$
- $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  : Training examples
- $\hat{f}(\mathbf{x})$  : Learned function
- $\mathcal{M}$  : Model

# Generalization Error

- We want to obtain  $\hat{f}(x)$  such that output values at unlearned test input points  $t$  can be accurately estimated.

- Suppose  $t \sim q(x)$

- Expected test error (generalization error):

$$G = \int_{\mathcal{D}} \left( \hat{f}(t) - f(t) \right)^2 q(t) dt$$

- **Goal:** Obtain  $\hat{f}(x)$  such that  $G$  is minimized.

# Formal Description of Problems<sup>6</sup>

$$G = \int_{\mathcal{D}} \left( \hat{f}(\mathbf{t}) - f(\mathbf{t}) \right)^2 q(\mathbf{t}) d\mathbf{t}$$

■ Active learning:  $\min_{\{\mathbf{x}_i\}_{i=1}^n} G$

■ Model selection:  $\min_{\mathcal{M}} G$

■ Learning methods:  $\min_{\hat{f} \in \mathcal{M}} G$

# Today's Plan

## ■ Models

- Linear models
- Kernel models

## ■ Learning methods

- Least-squares learning

# Linear/Non-Linear Models

- Model is a set of functions from which learning result functions are searched.
- We use a family of functions  $\hat{f}(x)$  parameterized by

$$\alpha = (\alpha_1, \alpha_2, \dots, \alpha_b)^\top$$

- **Linear model:**  $\hat{f}(x)$  is linear with respect to  $\alpha$   
(Note: **not necessarily linear with respect to  $x$** )
- **Non-linear model:** Otherwise



# Linear Models

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^b \alpha_i \varphi_i(\mathbf{x})$$

■  $\{\varphi_i(\mathbf{x})\}_{i=1}^b$  : Linearly independent functions

■ For example, when  $d = 1$

- Polynomial

$$1, x, x^2, \dots, x^{b-1}$$

- Trigonometric polynomial

$$1, \sin x, \cos x, \dots, \sin kx, \cos kx$$

$$b = 2k + 1$$

# Multi-Dimensional Linear Models<sup>10</sup>

- For multidimensional input ( $d > 1$ ), tensor product could be used.

$$\hat{f}(\mathbf{x}) = \sum_{i_1=1}^c \sum_{i_2=1}^c \cdots \sum_{i_d=1}^c$$

$$\alpha_{i_1, i_2, \dots, i_d} \varphi_{i_1}(x^{(1)}) \varphi_{i_2}(x^{(2)}) \cdots \varphi_{i_d}(x^{(d)})$$

$$\mathbf{x} = (x^{(1)}, x^{(2)}, \dots, x^{(d)})^\top$$

- The number of parameters is  $b = c^d$ , which increases exponentially w.r.t.  $d$ .
- Infeasible for large  $d$  !

# Additive Models

- For large  $d$ , we have to reduce the number of parameters.

- Additive model:

$$\hat{f}(\mathbf{x}) = \sum_{j=1}^d \sum_{i=1}^c \alpha_{i,j} \varphi_i(x^{(j)})$$

- The number of parameters is only  $b = cd$ .
- However, additive model is too simple so its representation capability may not be rich enough in some application.

# Kernel Models

## ■ Linear model:

$\{\varphi_i(\mathbf{x})\}_{i=1}^b$  do not depend on  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$

## ■ Kernel model:

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

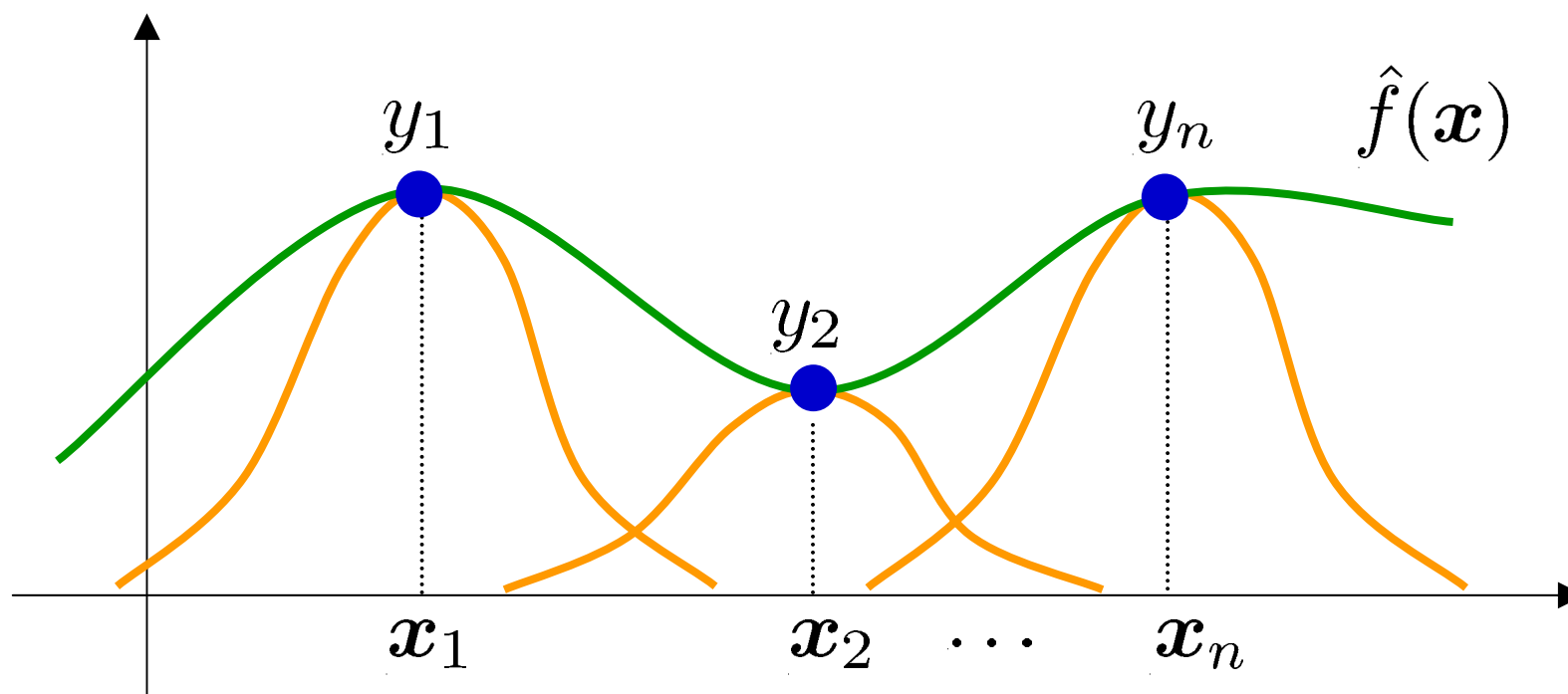
## ■ $K(\mathbf{x}, \mathbf{x}')$ : Kernel function

e.g., Gaussian kernel

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2c^2}\right)$$

# Kernel Models (cont.)

- Put kernel functions at training input points.



# Kernel Models (cont.)

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

- The number of parameters is  $n$  , which is **independent** of the input dimensionality  $d$  .
- Although kernel model is linear w.r.t.  $\alpha$  , the number of parameters grows as the number of training samples increases.
- Mathematical treatment could be different from ordinary linear models (called “non-parametric models” in statistics).

# Summary of Linear Models

- **Linear model (tensor):**  
High flexibility, high complexity
- **Linear model (additive):**  
Low flexibility, low complexity
- **Kernel model:**  
Moderate flexibility, moderate complexity
- Good model depends on applications.
- Later in model selection, we discuss how to choose appropriate models.

# Learning Methods

## ■ Linear learning methods:

Parameter vector  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_b)^\top$   
is estimated linearly w.r.t.

$$\mathbf{y} = (y_1, y_2, \dots, y_n)^\top$$

## ■ Non-linear learning methods: Otherwise



# Linear Learning for Linear Models / Kernel Models

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^b \alpha_i \varphi_i(\mathbf{x})$$

- In linear learning methods, a learned parameter vector is given by

$$\hat{\alpha} = L\mathbf{y}$$

$L$  : Learning matrix

# Least-Squares Learning

- Learn  $\alpha$  such that the squared error at training input points is minimized:

$$\hat{\alpha}_{LS} = \operatorname{argmin}_{\alpha \in \mathbb{R}^b} J_{LS}(\alpha)$$

$$\begin{aligned} J_{LS}(\alpha) &= \sum_{i=1}^n \left( \hat{f}(x_i) - y_i \right)^2 \\ &= \| \mathbf{X} \alpha - \mathbf{y} \|^2 \end{aligned}$$


$$\mathbf{X}_{i,j} = \varphi_j(x_i) : \text{Design matrix } (n \times b)$$

- In the following, we assume  $\operatorname{rank}(\mathbf{X}) = b$

# How to Obtain Solutions

## ■ Extreme-value condition:

$$\nabla J_{LS}(\hat{\alpha}_{LS}) = 2\mathbf{X}^\top (\mathbf{X}\hat{\alpha}_{LS} - \mathbf{y}) = 0$$


$$\hat{\alpha}_{LS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

## ■ Therefore, LS is linear learning.

$$\hat{\alpha}_{LS} = \mathbf{L}_{LS} \mathbf{y}$$

$$\mathbf{L}_{LS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

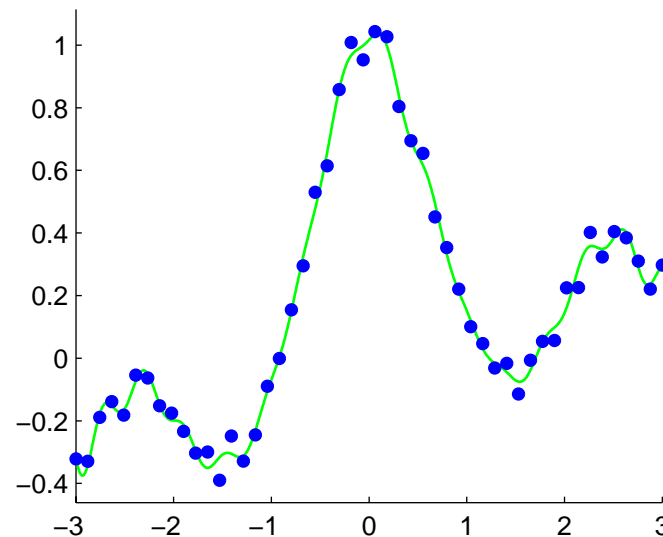
If you are not familiar with vector-derivatives, see e.g.,  
“Matrix Cookbook” (<http://2302.dk/uni/matrixcookbook.html>)

# Example of LS

$$\hat{f}(x) = \sum_{i=1}^b \alpha_i \varphi_i(x)$$

## ■ Trigonometric polynomial model

$1, \sin x, \cos x, \dots, \sin 15x, \cos 15x$  ( $b = 31$ )



# Homework

$$\hat{f}(\boldsymbol{x}) = \sum_{i=1}^n \alpha_i K(\boldsymbol{x}, \boldsymbol{x}_i)$$

1. Prove that the LS solution in kernel models is given by

$$\hat{\alpha}_{LS} = \boldsymbol{L}_{LS} \boldsymbol{y}$$

$$\boldsymbol{L}_{LS} = \boldsymbol{K}^{-1}$$

$$\boldsymbol{K}_{i,j} = K(\boldsymbol{x}_i, \boldsymbol{x}_j)$$

(Kernel matrix)

## Homework (cont.)

2. For your own toy 1-dimensional data, perform simulations using

- Gaussian kernel models
- Least-squares learning

and analyze the results when, e.g.,

- Target functions
- Number of samples
- Noise level
- Width of Gaussian kernel

are changed.

■ **Deadline: beginning of next class**