

# Least Squares

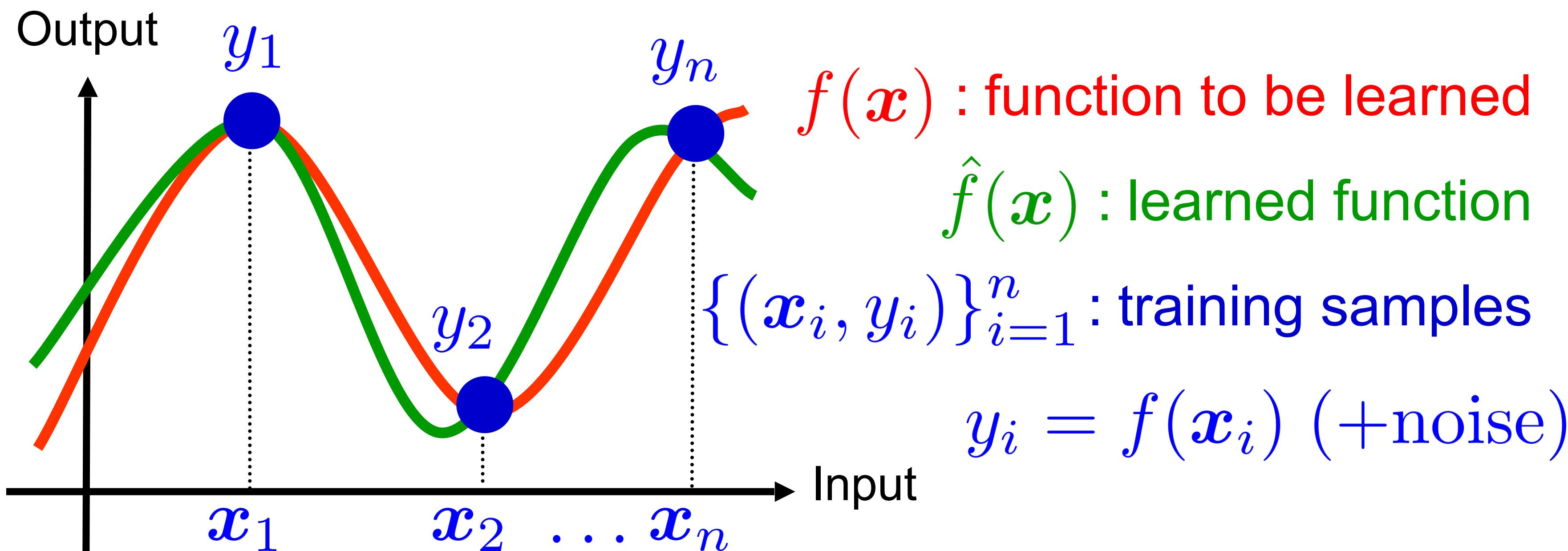
Masashi Sugiyama, Takashi Ishida  
[sugi@k.u-tokyo.ac.jp](mailto:sugi@k.u-tokyo.ac.jp), [ishi@k.u-tokyo.ac.jp](mailto:ishi@k.u-tokyo.ac.jp)  
<http://www.ms.k.u-tokyo.ac.jp>

# Schedule

- 1. 04/8 Introduction
- 2. 04/15 Regression 1
- 3. 04/22 Regression 2
- 04/30 Cancelled
- 4. 05/13 Classification 1
- 5. 05/20 Classification 2
- 6. 05/27 Deep learning 1
- 06/03 No lecture
- 7. 06/10 Deep learning 2
- 8. 06/17 Deep learning 3
- 9. 06/24 Semi-supervised learning
- 10. 07/01 Language models
- 11. 07/08 Representation learning 1
- 12. 07/15 Representation learning 2
- 13. 07/22 Advanced topics

# Regression = Function Approximation

3



Learn a function that is close to the  
underlying function by using training samples

# Contents

1. Learning model (Chap. 2)
  - A) Linear model
  - B) Kernel model
  - C) Non-linear model
2. Least squares regression (Chap. 3)
3. Regularized regression (Chap. 4)

# Linear/non-linear model

- **Model:** candidate set to search for a learned function
  - A model is specified when we choose parameter  $\theta$ 
$$\{f_{\theta}(x) \mid \theta = (\theta_1, \dots, \theta_b)^T\}$$
- **Linear model:**  $f_{\theta}(x)$  is linear w.r.t.  $\theta$   
( note: does not mean linear w.r.t. the input  $x$ )
- **Non-linear model:** other models

# Contents

6

1. learning model (Chap.2)
  - A) linear model
  - B) kernel model
  - C) non-linear model
2. least squares regression (Chap.3)
3. Regularized regression (Chap. 4)

# Linear model

- Usual form: 
$$f_{\theta}(x) = \sum_{j=1}^b \theta_j \phi_j(x)$$
 b functions
- $\{\phi_j(x)\}_{j=1}^b$ : known basis functions that are **linearly independent**
- Examples of basis functions when  $d = 1$ :
  - **Identity:**  $\phi(x) = x$  (model becomes linear in input)
  - **Polynomial functions:**  $(1, x, x^2, \dots, x^{b-1})$
  - **Sinusoidal functions:**  $(1, \sin x, \cos x, \dots, \sin kx, \cos kx)$

$$b = 2k + 1$$

# Multi-dimensional linear model

8

- Additive model:

$$\boldsymbol{x} = (x^{(1)}, \dots, x^{(d)})^T$$

T : transpose

 $d \geq 2$ Ref:  $x_i$  is i-th sample

$$f_{\theta}(\boldsymbol{x}) = \sum_{k=1}^d \sum_{j=1}^{b'} \theta_{k,j} \phi_j(x^{(k)})$$

- Number of parameters:  $b = b'd$

- Increases linearly w.r.t. input dimension  $d$ .
- Note: Multiplicative model (in the next slide) grows exponentially w.r.t.  $d$ .
- Since it only considers the sum of one-dimensional basis functions, it cannot represent complex functions with interactions.

# Multi-dimensional linear model

- Multiplicative model for multi-dimensional inputs:

$$f_{\theta}(\mathbf{x}) = \sum_{j_1=1}^{b'} \cdots \sum_{j_d=1}^{b'} \theta_{j_1, \dots, j_d} \phi_{j_1}(x^{(1)}) \cdots \phi_{j_d}(x^{(d)})$$

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

- # of Params:  $b = (b')^d$ 
  - Exponentially grows w.r.t. input dim  $d$
  - $b = 10^{100}$  when  $b' = 10, d = 100$
  - Not practical from a computational/statistical perspective
- Can only be used when  $d$  is small (We will discuss about sparsity in the sparse learning lecture!)

# Contents

10

1. Learning model (Chap.2)
  - A) Linear model
  - B) Kernel model
  - C) Non-linear model
2. Least squares regression (Chap.3)
3. Regularized regression (Chap. 4)

# Kernel model

11

- **Linear model** (that we discussed so far):

- Basis functions  $\{\phi_j(\mathbf{x})\}_{j=1}^b$  do not depend on the training samples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$

- **Kernel model**:

- We use training input samples  $\{\mathbf{x}_i\}_{i=1}^n$  to design the basis functions

$$f_{\theta}(\mathbf{x}) = \sum_{j=1}^n \theta_j K(\mathbf{x}, \mathbf{x}_j)$$

- Ex: **Gaussian kernel**

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

$$\|\mathbf{a}\|^2 = \mathbf{a}^\top \mathbf{a}$$

$h(>0)$ : bandwidth

# Kernel model

12

$$f_{\theta}(x) = \sum_{j=1}^n \theta_j K(x, x_j)$$

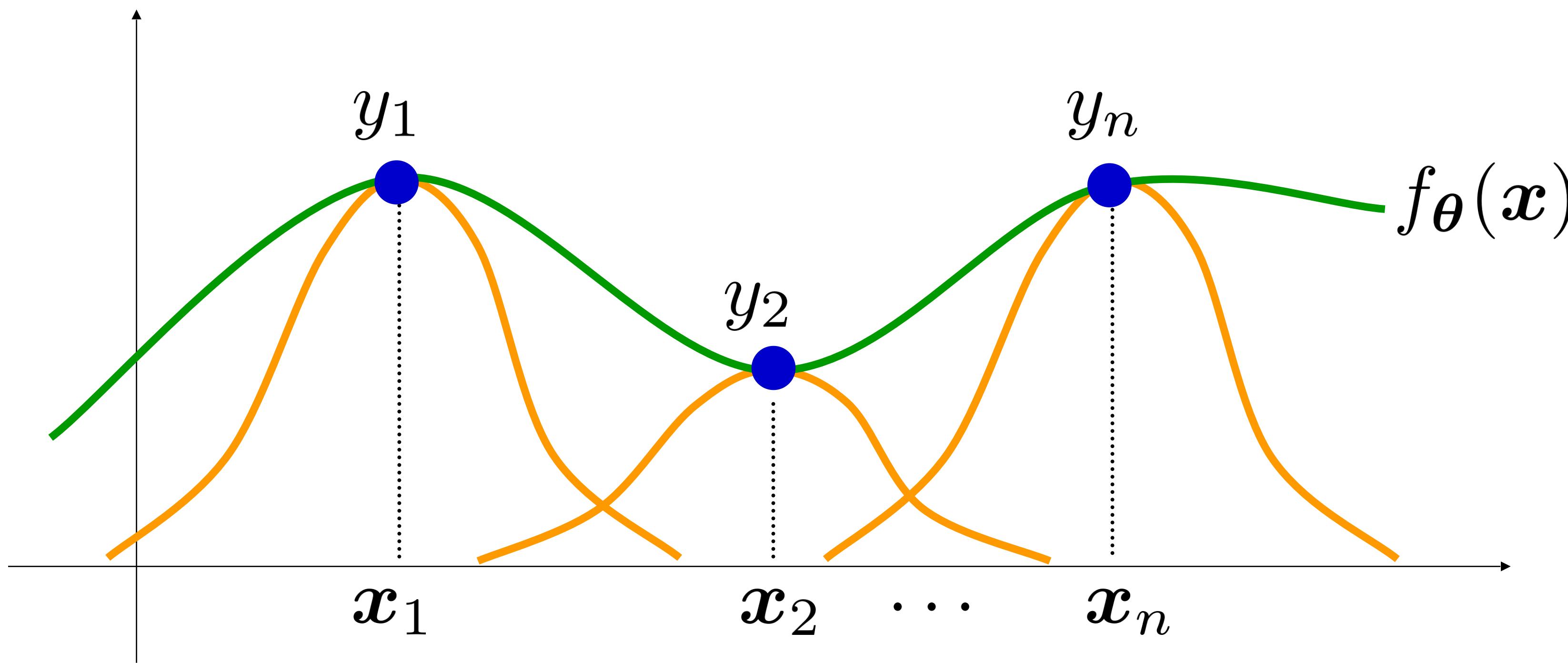
- The # of params is  $n$  and does not depend on input dim  $d$
- Linear in parameter:
  - kernel models is a type of linear model
- However, the # of params depend on the # of training samples, so have different statistical properties. Sometimes called **nonparametric models** due to this reason.
- In the scope of this course, you can treat kernel models as linear models.

# Gaussian kernel model

13

$$f_{\theta}(x) = \sum_{j=1}^n \theta_j \exp \left( -\frac{\|x - x_j\|^2}{2h^2} \right)$$

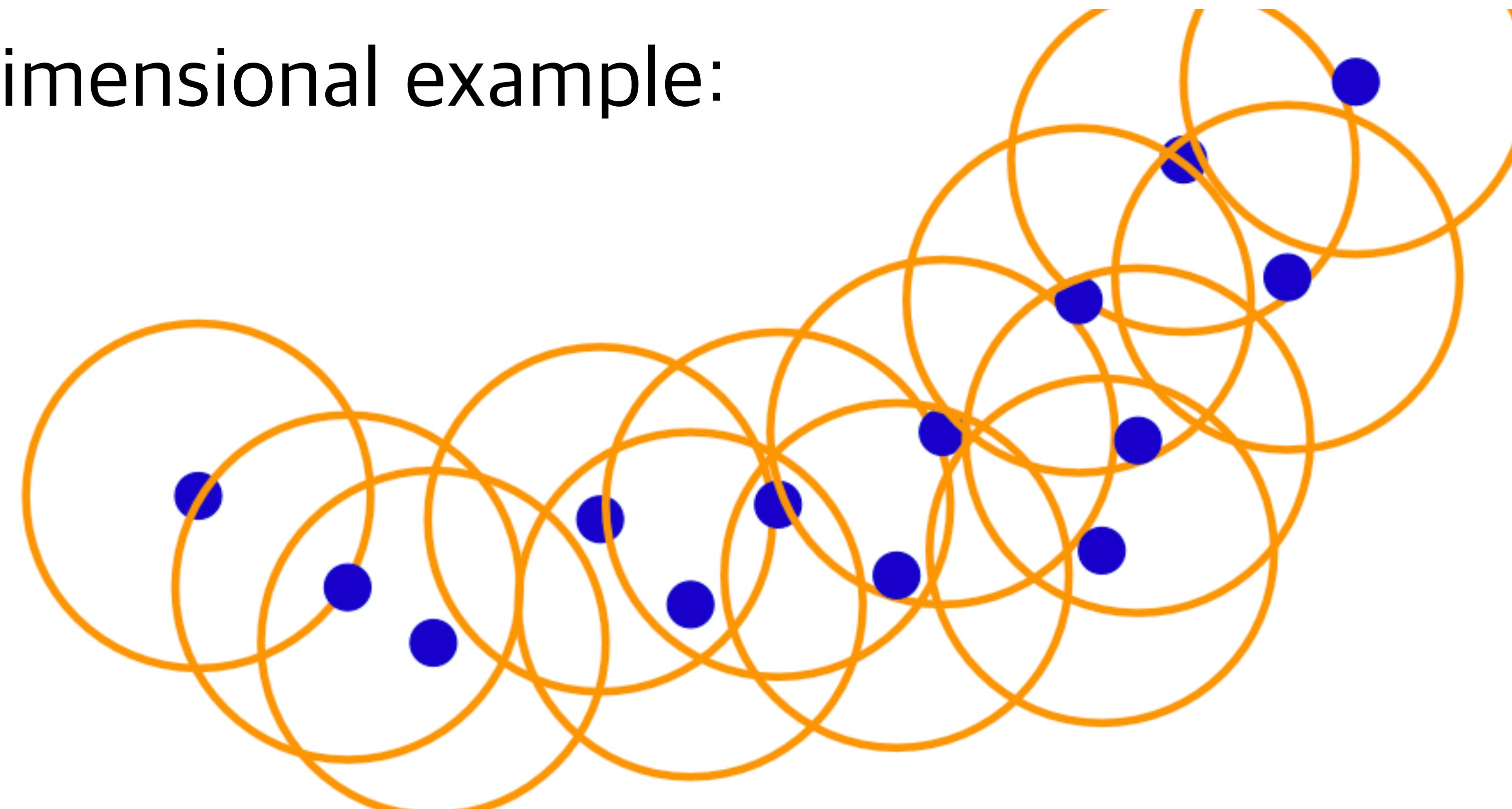
- Place Gaussian functions on each training input samples.



# Gaussian kernel model

14

- When the  $n$  training samples are biased in their distribution across the input space, the Gaussian kernel model automatically ignores regions where there are no training input samples.
- A 2-dimensional example:



# Contents

15

1. learning model (Chap.2)
  - A) linear model
  - B) kernel model
  - C) non-linear model
2. least squares regression (Chap.3)
3. Regularized regression (Chap. 4)

# Non-linear model

16

- Models that are not linear w.r.t. the parameters
  - An example: a linear model with the basis functions including the parameter.

$$f_{\theta}(\mathbf{x}) = \sum_{j=1}^b \alpha_j \phi(\mathbf{x}; \boldsymbol{\beta}_j) \quad \theta = (\boldsymbol{\alpha}^\top, \boldsymbol{\beta}_1^\top, \dots, \boldsymbol{\beta}_b^\top)^\top$$

$\phi(\mathbf{x}; \boldsymbol{\beta})$ : non-linear w.r.t.  $\boldsymbol{\beta}$

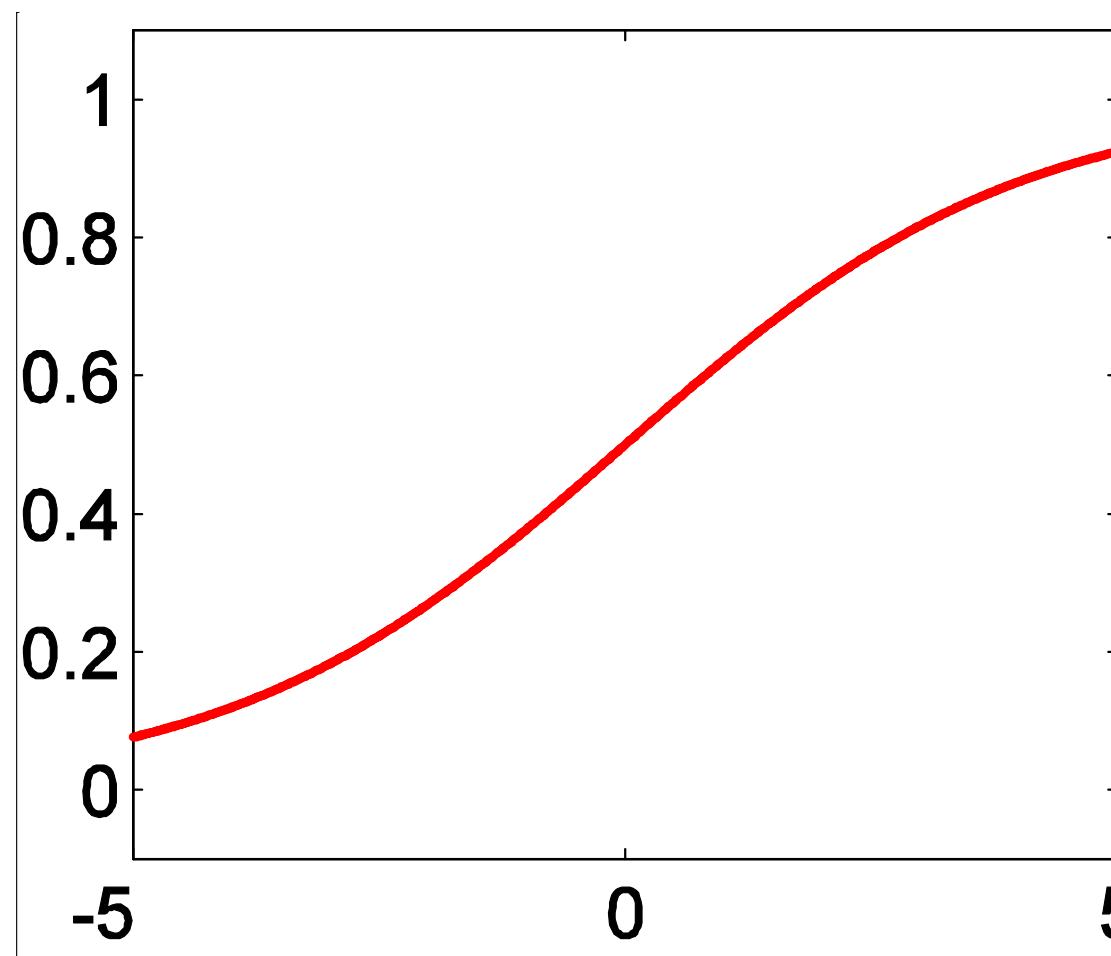
# Ex 1: basis function that includes param

17

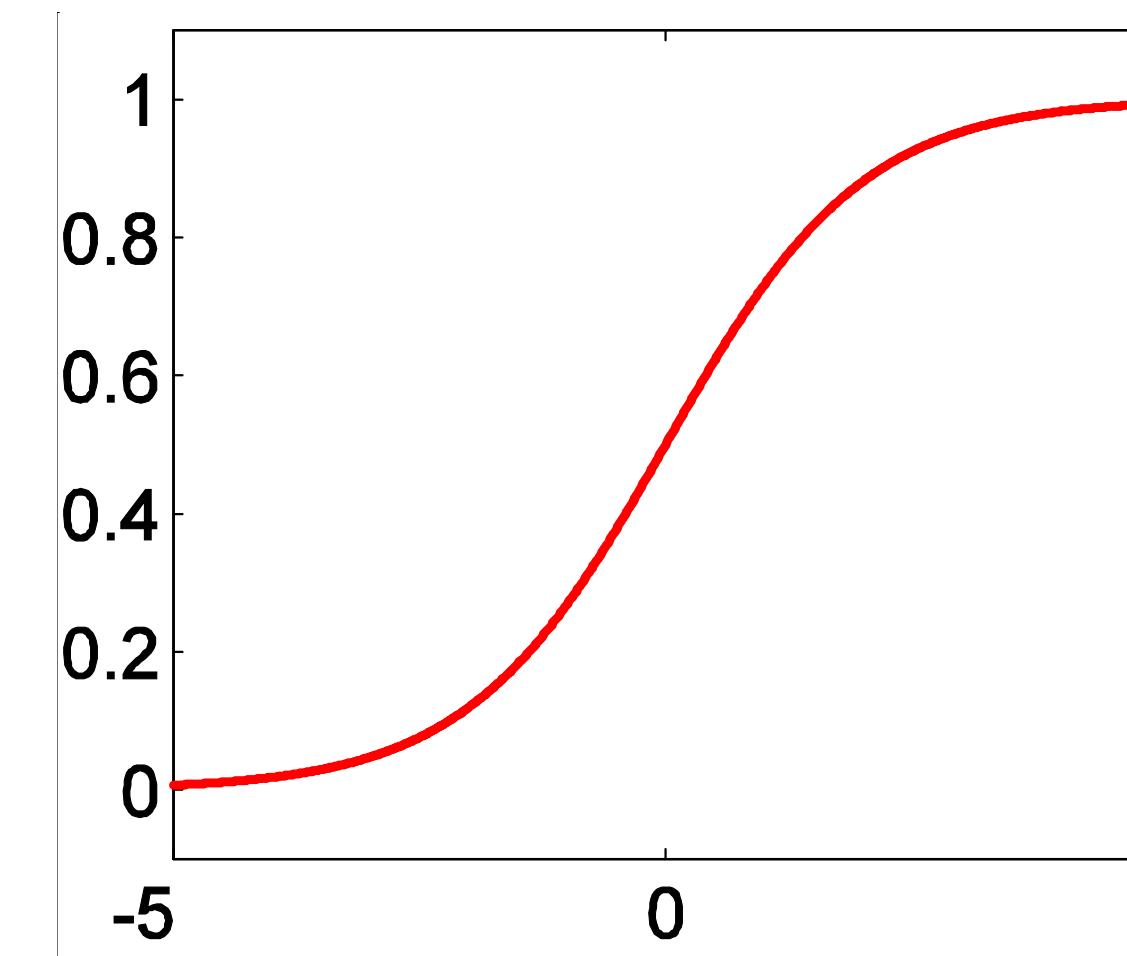
- Sigmoid function:

$$\phi(\mathbf{x}; \boldsymbol{\beta}) = \frac{1}{1 + \exp(-\mathbf{x}^\top \mathbf{w} - \gamma)} \quad \boldsymbol{\beta} = (\mathbf{w}^\top, \gamma)^\top$$

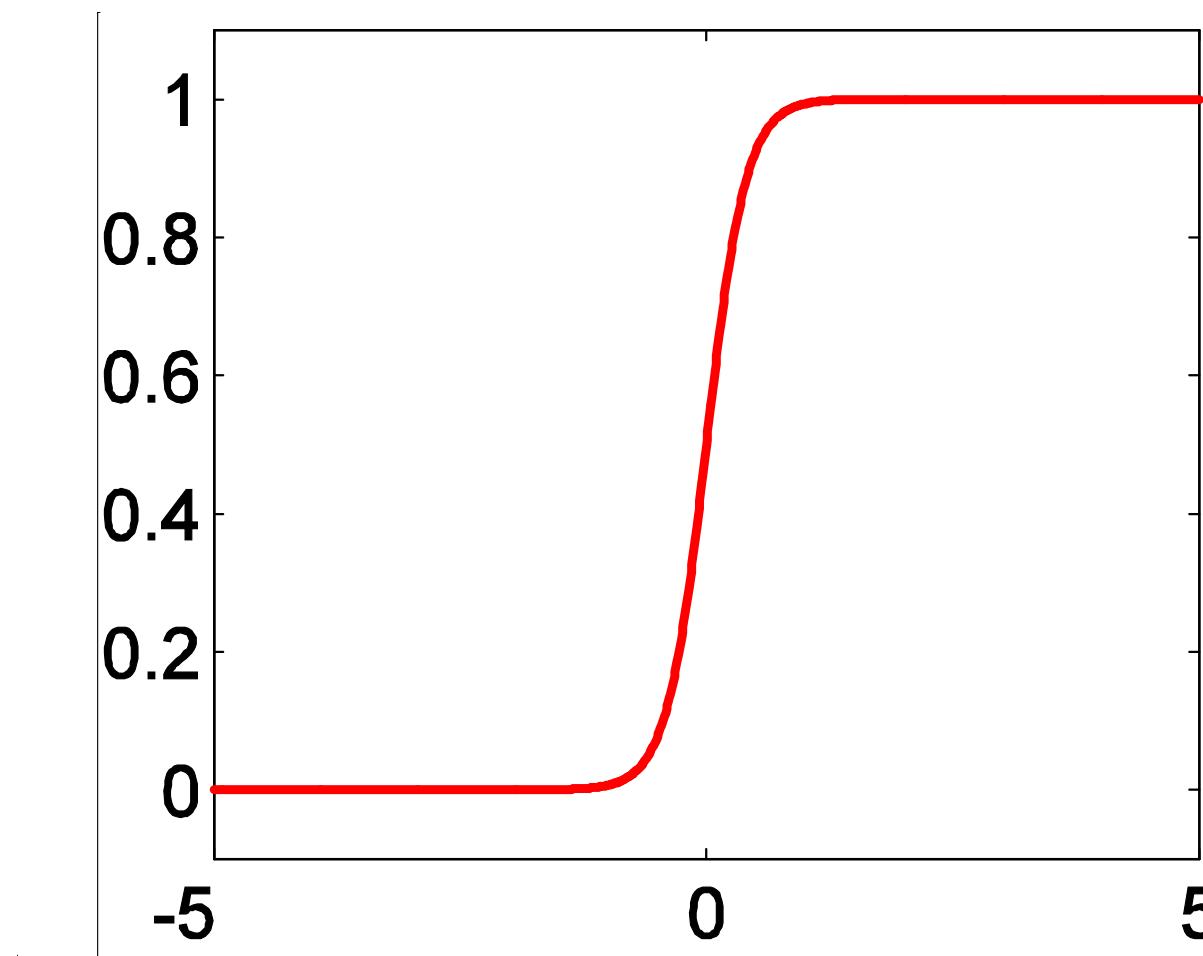
For  $d = 1$  and  $\gamma = 0$ :



$w$ : small



$w$ : medium



$w$ : large

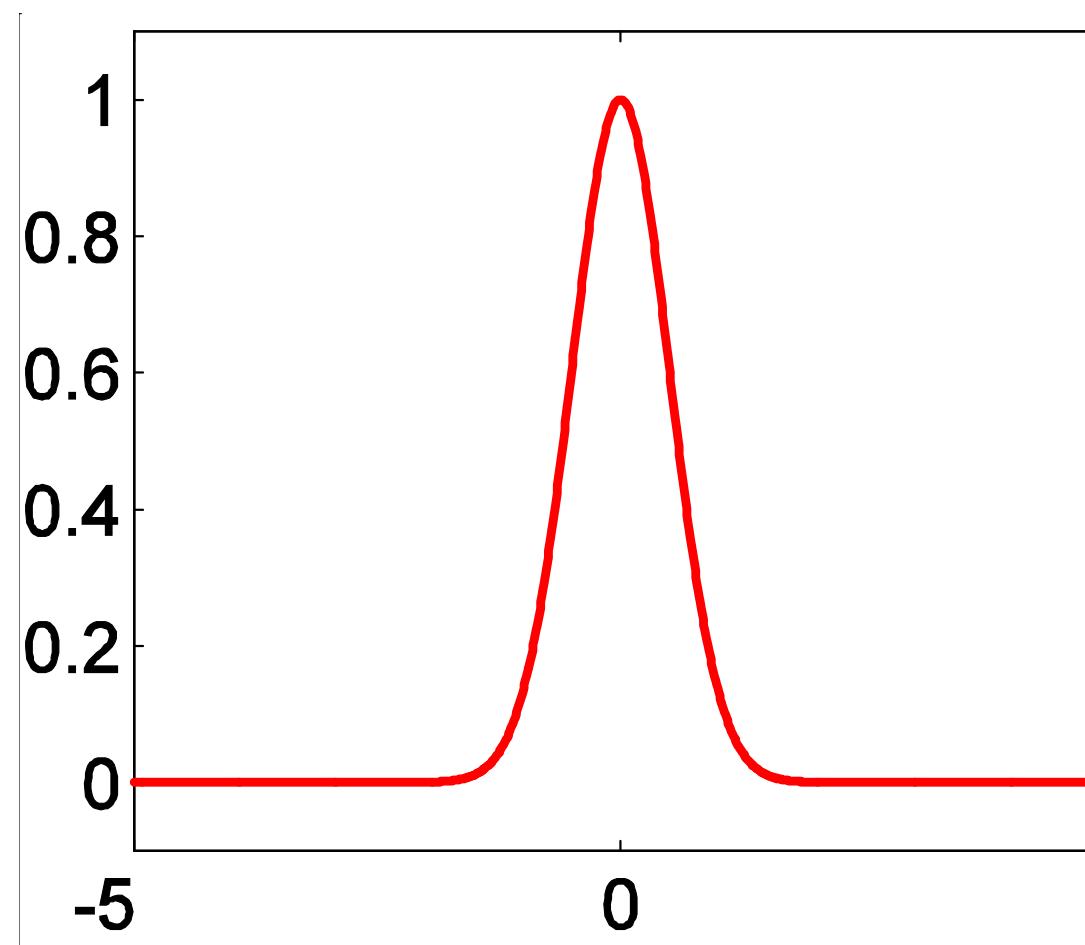
## Ex 2: basis function that includes param

18

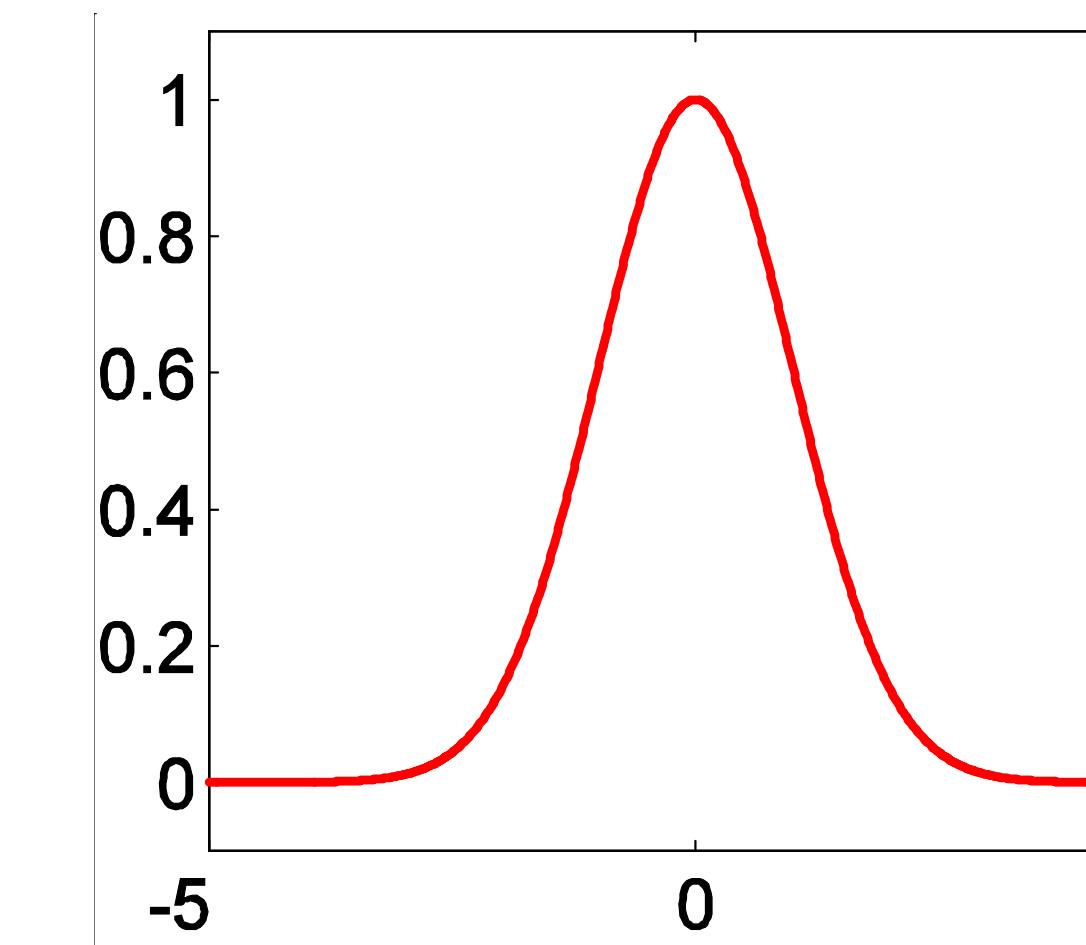
- **Gaussian kernel:** In contrast to linear Gaussian kernel models, we can learn the center and bandwidth of the Gaussian function

$$\phi(\mathbf{x}; \beta) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}\|^2}{2h^2}\right)$$

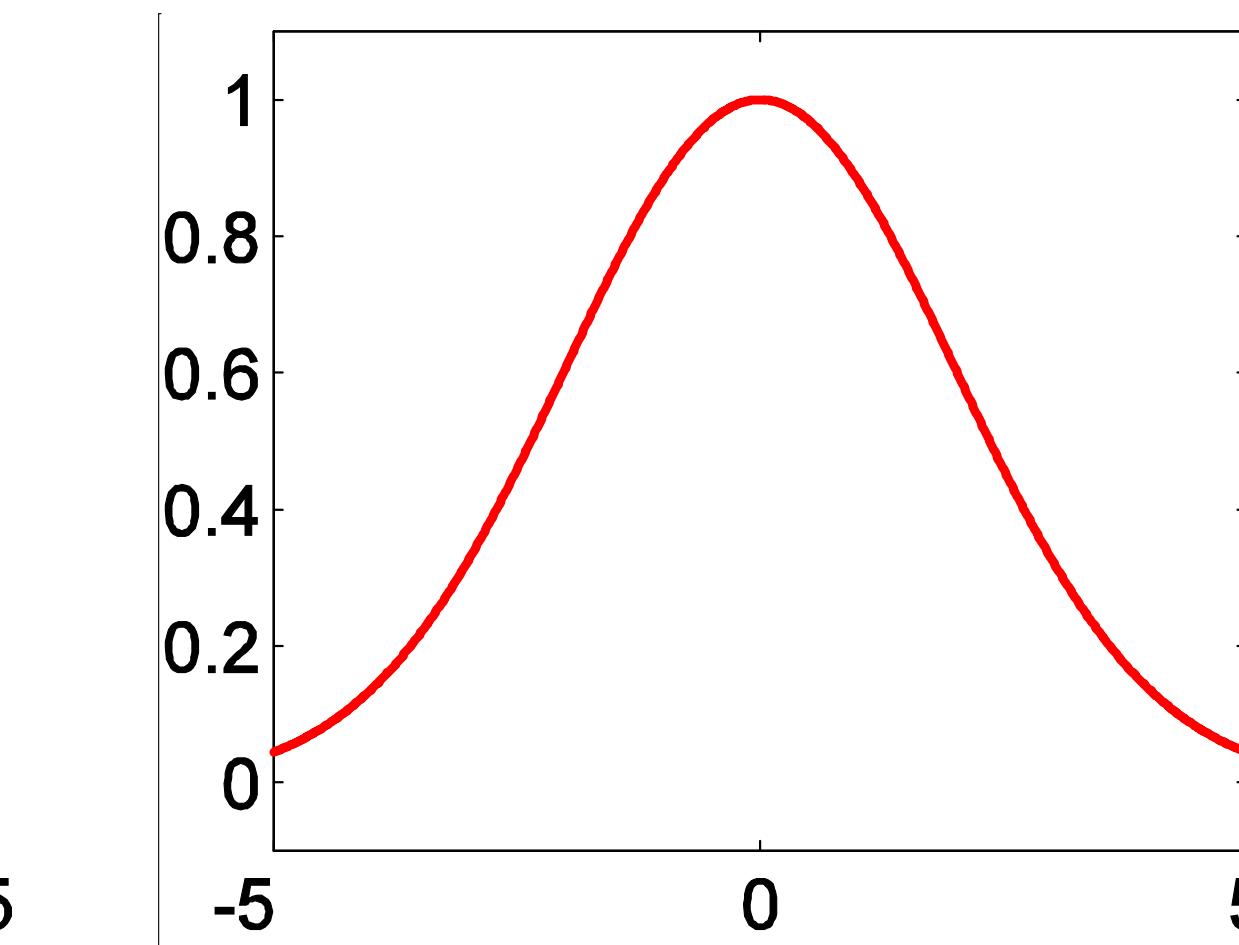
$$\beta = (\mathbf{c}^\top, h)^\top$$



$h$ : small



$h$ : medium



$h$ : large

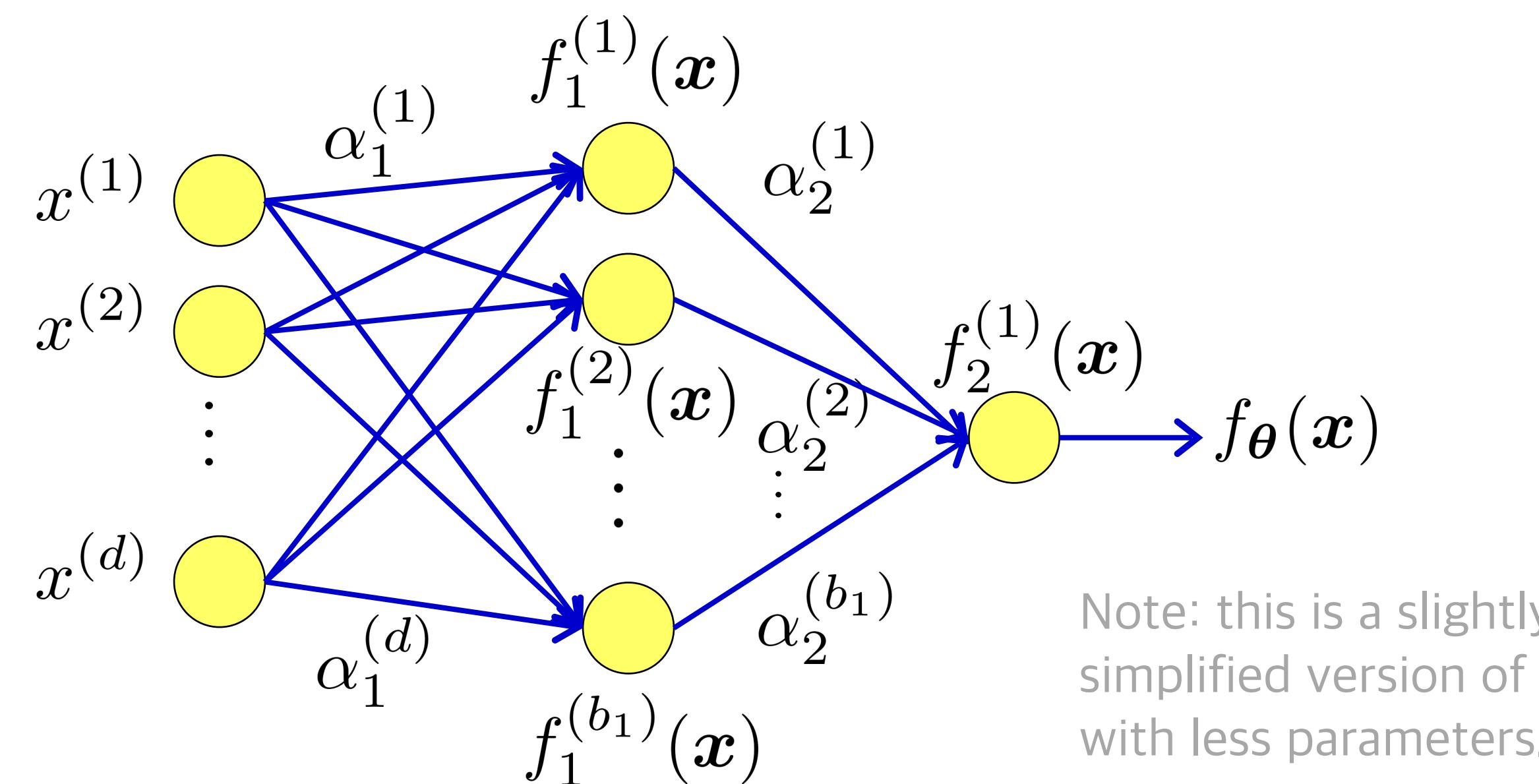
# Non-linear model

19

- Model that is non-linear w.r.t. the parameters
  - Hierarchical model

$$f_k^{(i)}(\mathbf{x}) = \sum_{j=1}^{b_{k-1}} \alpha_k^{(j)} \phi(f_{k-1}^{(j)}(\mathbf{x})), \quad f_0^{(i)}(\mathbf{x}) = x^{(i)}$$

$$f_{\theta}(\mathbf{x}) = f_K^{(1)}(\mathbf{x}) \quad \theta = (\alpha_1^{\top}, \dots, \alpha_K^{\top})^{\top}$$

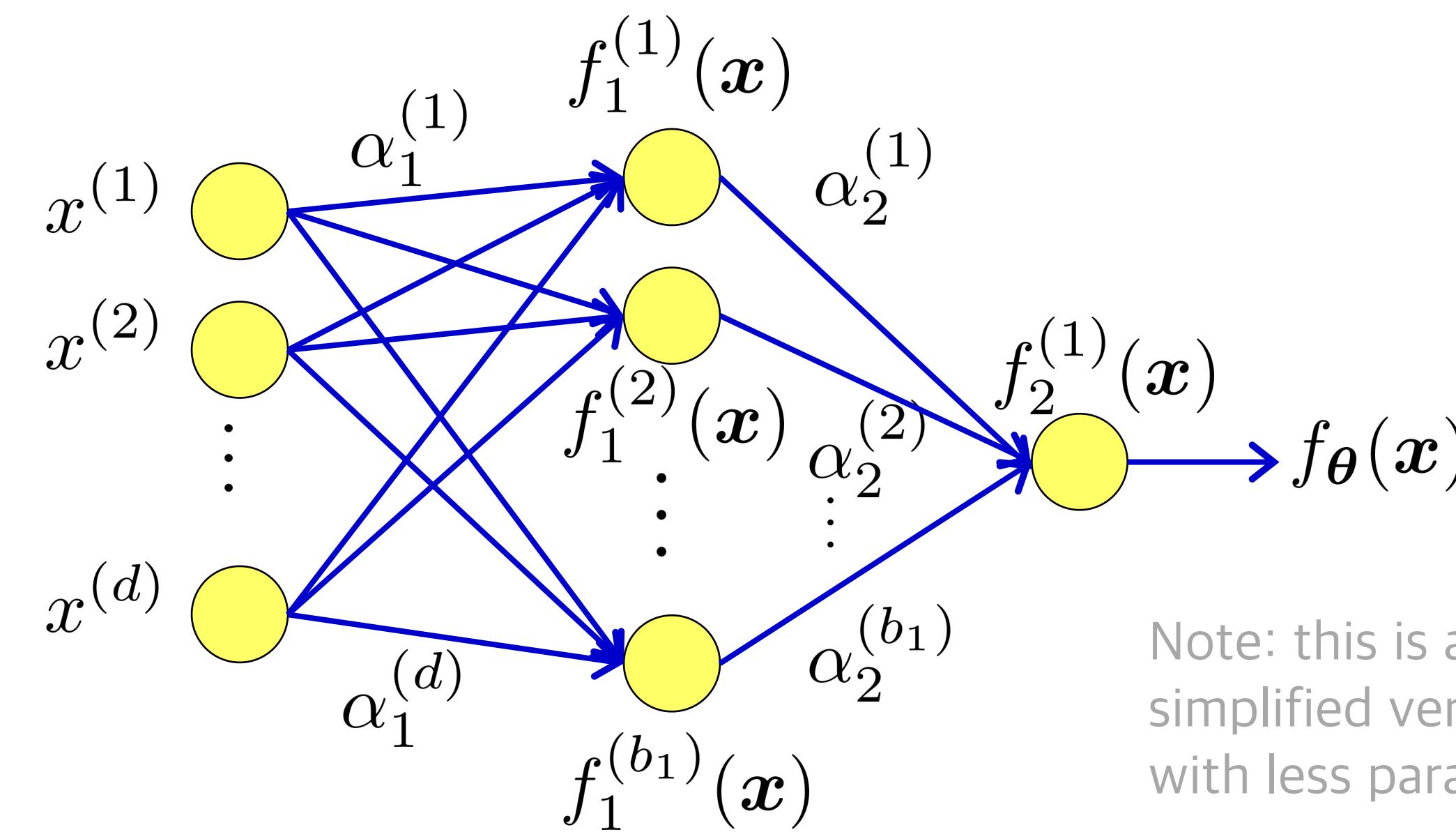


Note: this is a slightly simplified version of an MLP with less parameters.

# Some history on perceptrons

20

- The human brain consists of countless **nerve cells** connected in a network-like pattern.
- Sigmoid function behaves similarly to the **nerve cells**.
- Artificial neural networks are also called **perceptrons**.
- Mathematically, a 3-layer perceptron can **approximate any function**.



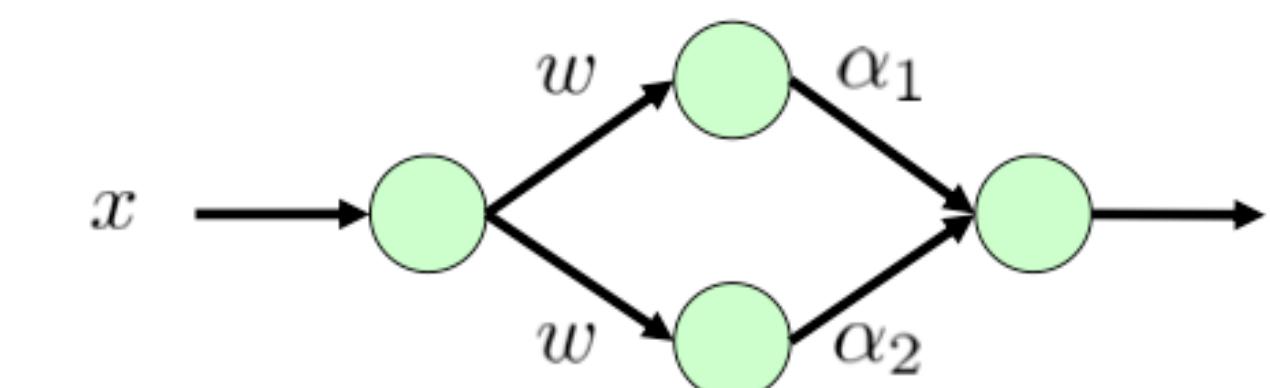
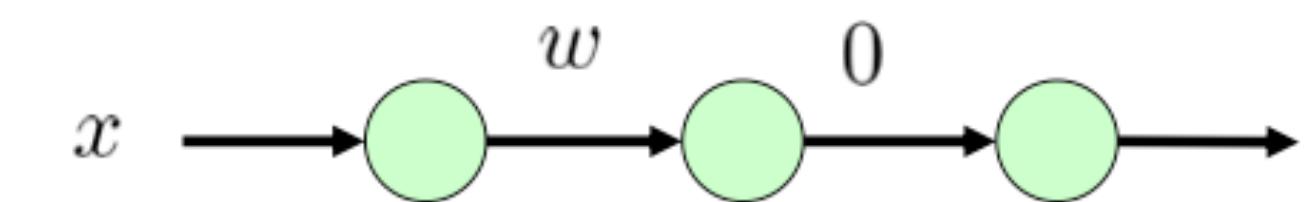
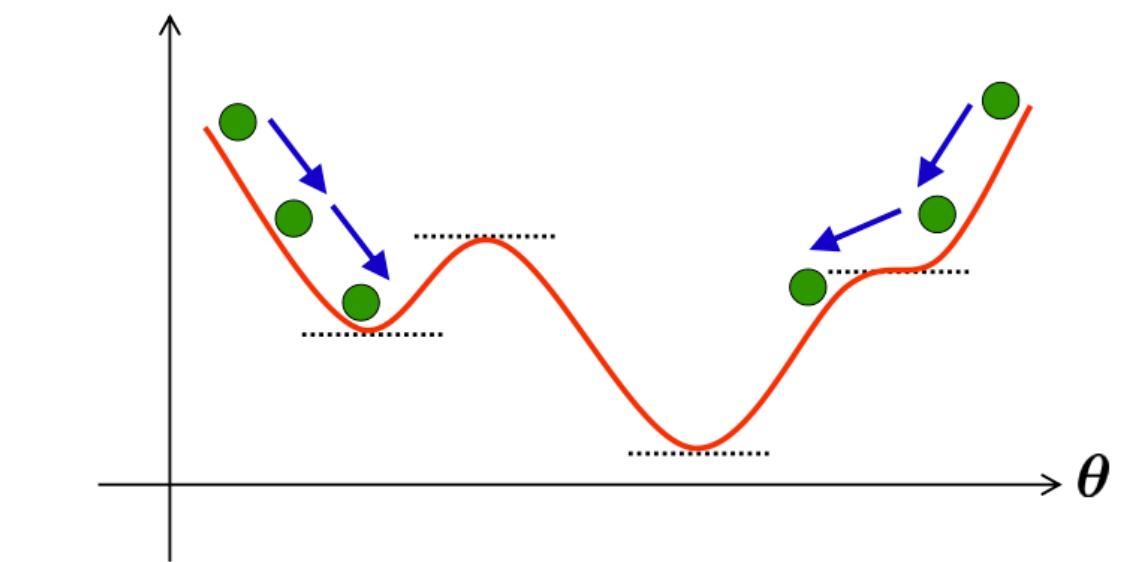
Note: this is a slightly simplified version of an MLP with less parameters.

# Difficulty of training hierarchical models

21

Historically, hierarchical models were **difficult to train!**

- Due to the existence of numerous local optima, finding the global optimum is difficult.
- Parameters and functions do not have a one-to-one correspondence.
  - When the weights of the second layer are zero, changing the weights of the first layer does not affect the function.
  - When the weights of the first layer are equal, if the sum of the weights of the second layer  $\alpha_1 + \alpha_2$  is a constant, the function remains the same.



Modern deep learning largely mitigate these training difficulties, and we will discuss this in later lectures.

# Summary of learning models

22

- Linear model: linear w.r.t. the parameters
  - Additive model: small number of params, low flexibility
  - Multiplicative model: High flexibility but too many params
  - kernel model: moderate flexibility and number of params
- Non-linear model: Non-linear w.r.t. to parameters
  - Hierarchical model: we will discuss this in the future
- We need to perform “model selection” since the decision will depend on the application.
- In this course: mainly use kernel models, occasionally hierarchical models such as neural network models.

# Contents

23

1. Learning model (Chap.2)
2. Least squares regression (Chap.3)
3. Regularized regression (Chap. 4)

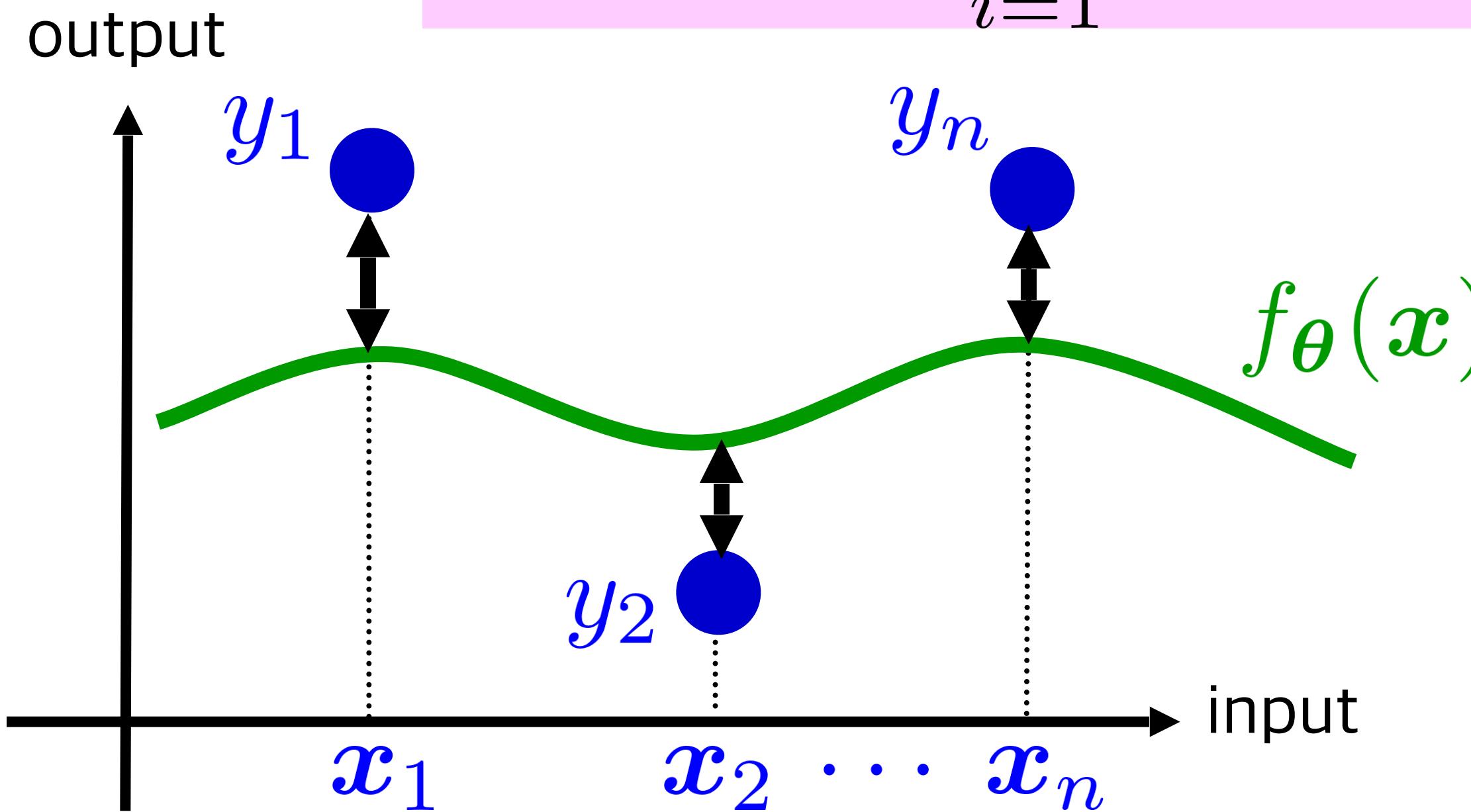
# Least squares regression

24

- **Objective:** minimize the squared error between the training outputs

$$\hat{\theta}_{\text{LS}} = \underset{\theta}{\operatorname{argmin}} J_{\text{LS}}(\theta)$$

$$J_{\text{LS}}(\theta) = \frac{1}{2} \sum_{i=1}^n \left( f_{\theta}(x_i) - y_i \right)^2$$



# Deriving the solution

- Rewrite the objective for model  $f_{\theta}(\mathbf{x}) = \sum_{j=1}^b \theta_j \phi_j(\mathbf{x})$

$$J_{\text{LS}}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n \left( f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i \right)^2$$

$$= \frac{1}{2} \|\Phi\boldsymbol{\theta} - \mathbf{y}\|^2 = \frac{1}{2} (\Phi\boldsymbol{\theta} - \mathbf{y})^\top (\Phi\boldsymbol{\theta} - \mathbf{y})$$

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

$$\Phi = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \cdots & \phi_b(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_1(\mathbf{x}_n) & \cdots & \phi_b(\mathbf{x}_n) \end{pmatrix} : \text{Design matrix}$$

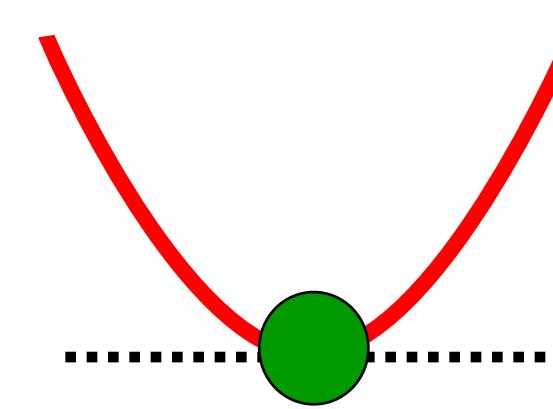
# Deriving the solution (2)

- Partial derivative is zero:

$$\nabla_{\theta} J_{\text{LS}} = \left( \frac{\partial J_{\text{LS}}}{\partial \theta_1}, \dots, \frac{\partial J_{\text{LS}}}{\partial \theta_b} \right)^{\top} = \Phi^{\top} \Phi \theta - \Phi^{\top} y = 0$$

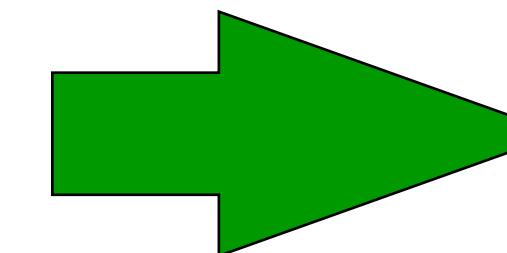
$$\frac{\partial}{\partial \theta} \theta^{\top} A \theta = 2A\theta$$

$$\frac{\partial}{\partial \theta} b^{\top} \theta = b$$



- Equation satisfied by the solution:

$$\Phi^{\top} \Phi \theta = \Phi^{\top} y$$



$$\hat{\theta}_{\text{LS}} = (\Phi^{\top} \Phi)^{-1} \Phi^{\top} y$$

Note: assume  $(\Phi^{\top} \Phi)^{-1}$  exists

- We analytically obtained the solution!

# Why is $\Phi^T\Phi$ non-singular?

- Recall that  $\{\phi_j(\mathbf{x})\}_{j=1}^b$  are known basis functions that are **linearly independent**. So the rank of  $\Phi$  is  $b$ .
- We first derive that  $\Phi^T\Phi$  is positive definite:
  - If  $\mathbf{a} \equiv \Phi\mathbf{c}$ , then  $\mathbf{c}^T\Phi^T\Phi\mathbf{c} = \mathbf{a}^T\mathbf{a} = \sum_{i=1}^n a_i^2 \geq 0$ .
  - Since  $\Phi$  is full column rank,  $\mathbf{a} \neq \mathbf{0}$  for any  $\mathbf{c} \neq \mathbf{0}$ .
  - Therefore we have  $\mathbf{c}^T\Phi^T\Phi\mathbf{c} > 0$ .
  - Since  $\Phi^T\Phi$  is positive definite, it is **non-singular**.
    - Proof by contradiction: If  $\Phi^T\Phi$  is singular, it means there exists a  $\mathbf{c} \neq \mathbf{0}$  such that  $\Phi^T\Phi\mathbf{c} = \mathbf{0}$ . Then  $\mathbf{c}^T\Phi^T\Phi\mathbf{c} = \mathbf{c}^T\mathbf{0} = 0$ .

# Second-order condition

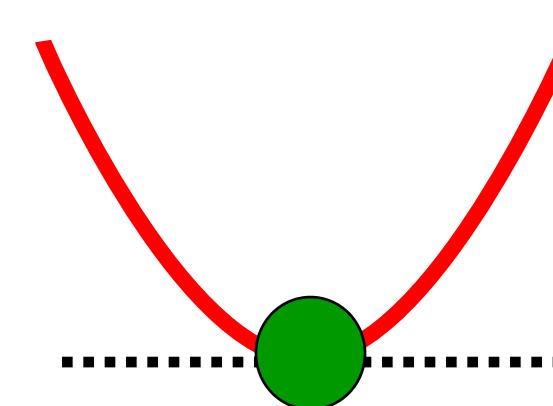
- To make sure the solution achieves a minimum instead of the maximum, we can check the second-order condition.
- We already derived:

$$\nabla_{\theta} J_{LS} = \Phi^T \Phi \theta - \Phi^T y$$

- We can further derive the Hessian:

$$\nabla_{\theta}^2 J_{LS} = \Phi^T \Phi$$

- From the previous slide, we know the Hessian is **positive definite**.



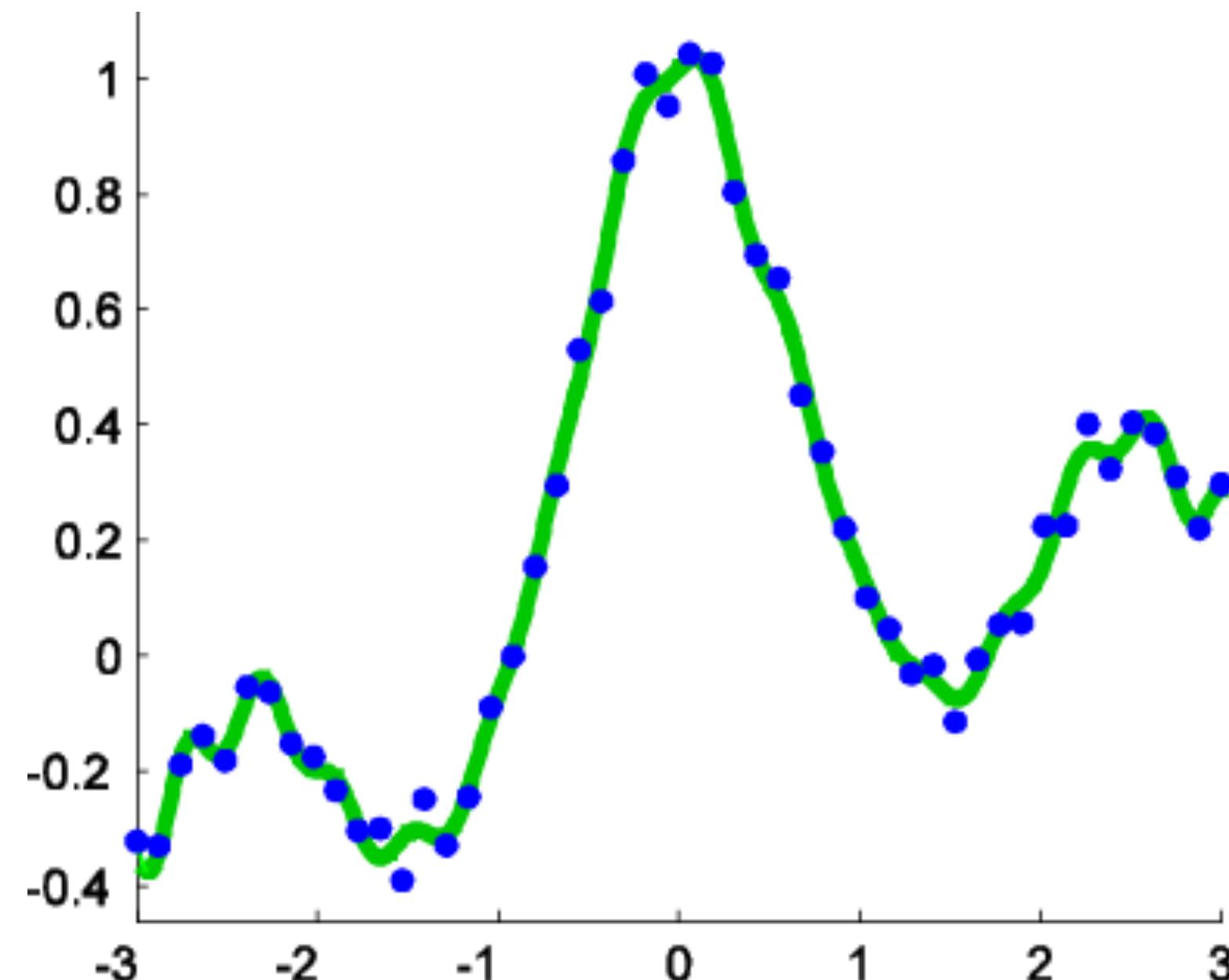
# Example

29

- An example of least squares regression with a sinusoidal model  $f_{\theta}(x) = \sum_{j=1}^b \theta_j \phi_j(x)$

$$\phi_i(x) = 1, \sin(x/2), \cos(x/2), \dots, \sin(15x/2), \cos(15x/2)$$

true function:  $f(x) = \sin(\pi x)/(\pi x) + 0.1x$



# Example with Python

```
import numpy as np; import matplotlib; matplotlib.use('TkAgg'); import
matplotlib.pyplot as plt; np.random.seed(0) # set seed for reproducibility

def generate_sample(xmin, xmax, sample_size):
    x = np.linspace(start=xmin, stop=xmax, num=sample_size)
    pix = np.pi * x
    target = np.sin(pix) / pix + 0.1 * x
    noise = 0.05 * np.random.normal(loc=0., scale=1., size=sample_size)
    return x, target + noise

def calc_design_matrix(x):
    sample_size = len(x)
    phi = np.empty(shape=(sample_size, 31)) # design matrix
    phi[:, 0] = 1.
    phi[:, 1::2] = np.sin(x[:, None] * np.arange(1, 16)[None] / 2)
    phi[:, 2::2] = np.cos(x[:, None] * np.arange(1, 16)[None] / 2)
    return phi

# create sample
sample_size = 50; xmin, xmax = -3, 3
x, y = generate_sample(xmin=xmin, xmax=xmax, sample_size=sample_size)

phi = calc_design_matrix(x) # calculate design matrix

theta = np.linalg.solve(phi.T.dot(phi), phi.T.dot(y[:, None])) # solve LS

# create data to visualize the prediction
X = np.linspace(start=xmin, stop=xmax, num=5000)
Phi = calc_design_matrix(X)
prediction = Phi.dot(theta)

# visualization
plt.clf(); plt.scatter(x, y, c='green', marker='o')
plt.plot(X, prediction); plt.show()
```

# Exercise

- Prove that when the noise in the training output follows a normal distribution independently, the least squares regression is consistent with the maximum likelihood estimation of the Gaussian model.

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^n \quad y_i = f(\mathbf{x}_i) + \epsilon_i \quad \epsilon_i \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2)$$

- Probabilistic model:

$$p_{\boldsymbol{\theta}}(y|\mathbf{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - f_{\boldsymbol{\theta}}(\mathbf{x}))^2}{2\sigma^2}\right)$$

- Log-likelihood:

$$L(\boldsymbol{\theta}) = \sum_{i=1}^n \log p_{\boldsymbol{\theta}}(y_i|\mathbf{x}_i)$$

- Likelihood estimation:

$$\max_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

$$J_{\text{LS}}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i)^2$$

# Justification of least squares regression

- When the noise in the output follows a normal distribution independently:
  - **Consistency:** the least squares solution converges in probability to the optimal solution.
  - If there are an infinite number of samples, the optimal parameters can be obtained.
- **Asymptotic efficiency:** has the smallest asymptotic variance among asymptotically normal estimators.
- When there are a sufficiently large number of samples, the estimation results are stable.

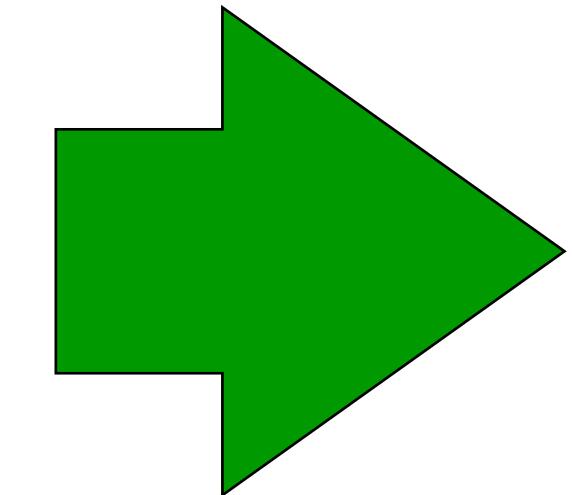
# Least squares regression w/ kernel models

$$f_{\theta}(x) = \sum_{j=1}^n \theta_j K(x, x_j)$$

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

- Least squares objective:

$$J_{\text{LS}}(\theta) = \frac{1}{2} \|K\theta - y\|^2$$



$$\hat{\theta}_{\text{LS}} = K^{-1}y$$

Kernel matrix:

$$K = \begin{pmatrix} K(x_1, x_1) & \cdots & K(x_1, x_n) \\ \vdots & \ddots & \vdots \\ K(x_n, x_1) & \cdots & K(x_n, x_n) \end{pmatrix}$$

# Details of derivation

- Objective and solution:

$$J_{\text{LS}}(\theta) = \frac{1}{2} \|K\theta - y\|^2 \longrightarrow \hat{\theta}_{\text{LS}} = K^{-1}y$$

- Since partial derivative is zero, we can derive:

$$\hat{\theta}_{\text{LS}} = (K^\top K)^{-1} K^\top y$$

- A property of Kernel matrix:  $K^\top = K$

- Hence,  $\hat{\theta}_{\text{LS}} = (K^2)^{-1} Ky = K^{-1}y$

$$K = \begin{pmatrix} K(x_1, x_1) & \cdots & K(x_1, x_n) \\ \vdots & \ddots & \vdots \\ K(x_n, x_1) & \cdots & K(x_n, x_n) \end{pmatrix}$$

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

# Summary of least squares regression

- **Main objective:** minimize the squared error between the training output and the predictions.
- **Regression:** basic technique for modeling the relationship between input and output variables.
- **When noise in the output is Gaussian:** the method can be interpreted as maximum likelihood estimation.
- **Analytical solutions:** for linear models and kernel models, we can find the optimal parameters analytically.

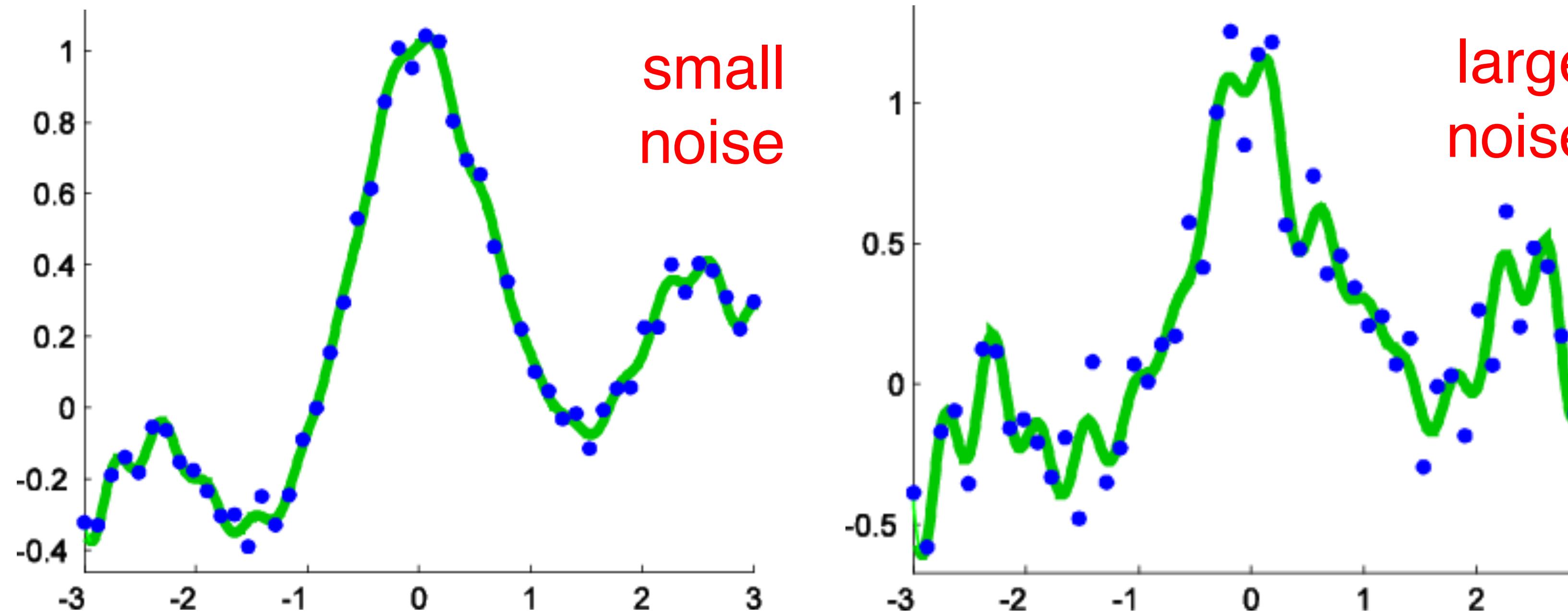
# Contents

36

1. Learning model (Chap.2)
2. Least squares regression (Chap.3)
3. Regularized regression (Chap. 4)
  - A) Constrained least squares regression
  - B) Model selection

# Issue of least squares regression

- **Overfitting:** occurs when a model adapts too much to the noise in the training samples.



$$\phi_i(x) = 1, \sin x/2, \cos x/2, \dots, \sin 15x/2, \cos 15x/2$$

→ properly constrain the model to prevent this

# Contents

38

1. Learning model (Chap.2)
2. Least squares regression (Chap.3)
3. Regularized regression (Chap. 4)
  - A) Constrained least squares regression
  - B) Model selection

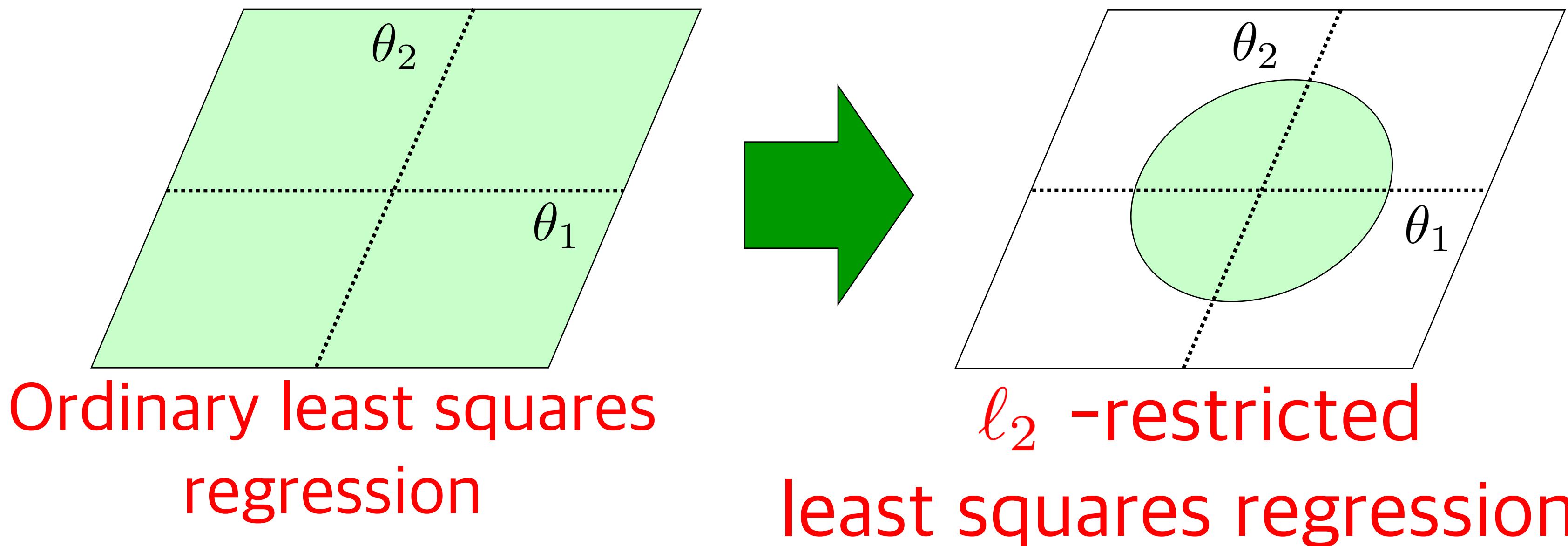
# $\ell_2$ –restricted least squares regression

39

- Avoid overfitting by restricting the model to be within a **hyper-cube**.  
(Because we want the absolute values of each elements of  $\theta$  to be small!)

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n (f_{\theta}(x_i) - y_i)^2 \text{ subject to } \|\theta\|^2 \leq R$$

$$R \geq 0 \quad f_{\theta}(x) = \sum_{j=1}^b \theta_j \phi_j(x)$$



# Deriving the solution

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n (f_{\theta}(x_i) - y_i)^2 \text{ subject to } \|\theta\|^2 \leq R$$

- Equivalent expression:

$$\min_{\theta} \left[ \frac{1}{2} \sum_{i=1}^n (f_{\theta}(x_i) - y_i)^2 + \frac{\lambda}{2} \|\theta\|^2 \right]$$

- $\lambda (> 0)$  : constant that is specified by  $R$
- We can use  $\lambda$  instead of  $R$  in practice.

# Interpretation

41

$$\min_{\theta} \left[ \frac{1}{2} \sum_{i=1}^n (f_{\theta}(x_i) - y_i)^2 + \frac{\lambda}{2} \|\theta\|^2 \right]$$



goodness of fit for training output

penalty term for preventing param values from becoming too large (regularization)

- The model strikes a good balance between fitting the training output well and keeping the parameter values small.
- Also called  $\ell_2$ -regularized regression

# Exercise

$$\hat{\theta} = \operatorname{argmin}_{\theta} \left[ \frac{1}{2} \sum_{i=1}^n \left( f_{\theta}(\mathbf{x}_i) - y_i \right)^2 + \frac{\lambda}{2} \|\theta\|^2 \right]$$

- Find the solution for linear model  $f_{\theta}(\mathbf{x}) = \sum_{j=1}^b \theta_j \phi_j(\mathbf{x})$

- Info:**  $\frac{1}{2} \sum_{i=1}^n \left( f_{\theta}(\mathbf{x}_i) - y_i \right)^2 = \frac{1}{2} \|\Phi\theta - \mathbf{y}\|^2$

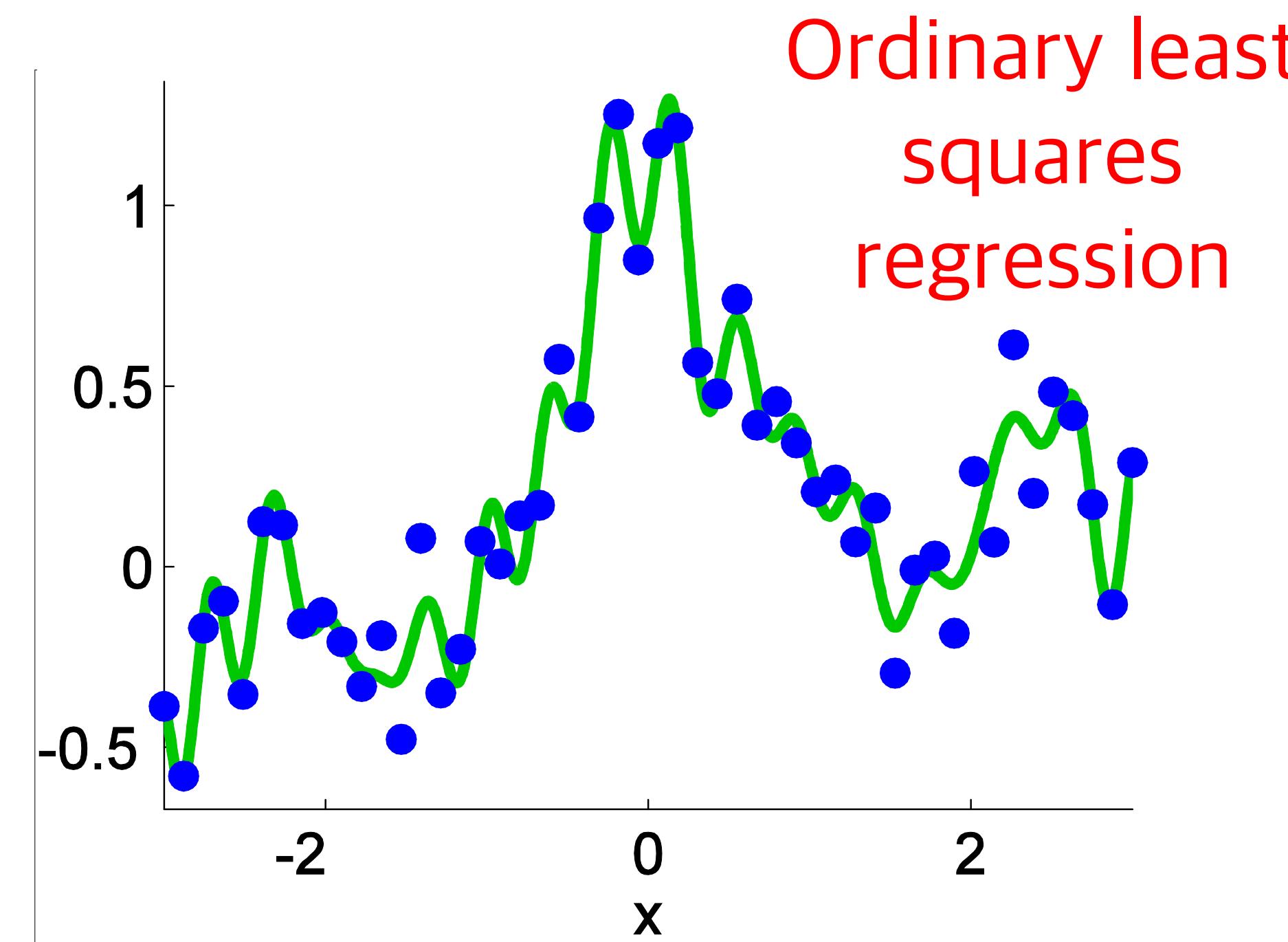
$$\Phi = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \cdots & \phi_b(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_1(\mathbf{x}_n) & \cdots & \phi_b(\mathbf{x}_n) \end{pmatrix} \quad \mathbf{y} = (y_1, \dots, y_n)^\top$$

# Example

- (Gaussian) kernel model

$$f_{\theta}(x) = \sum_{j=1}^n \theta_j K(x, x_j)$$

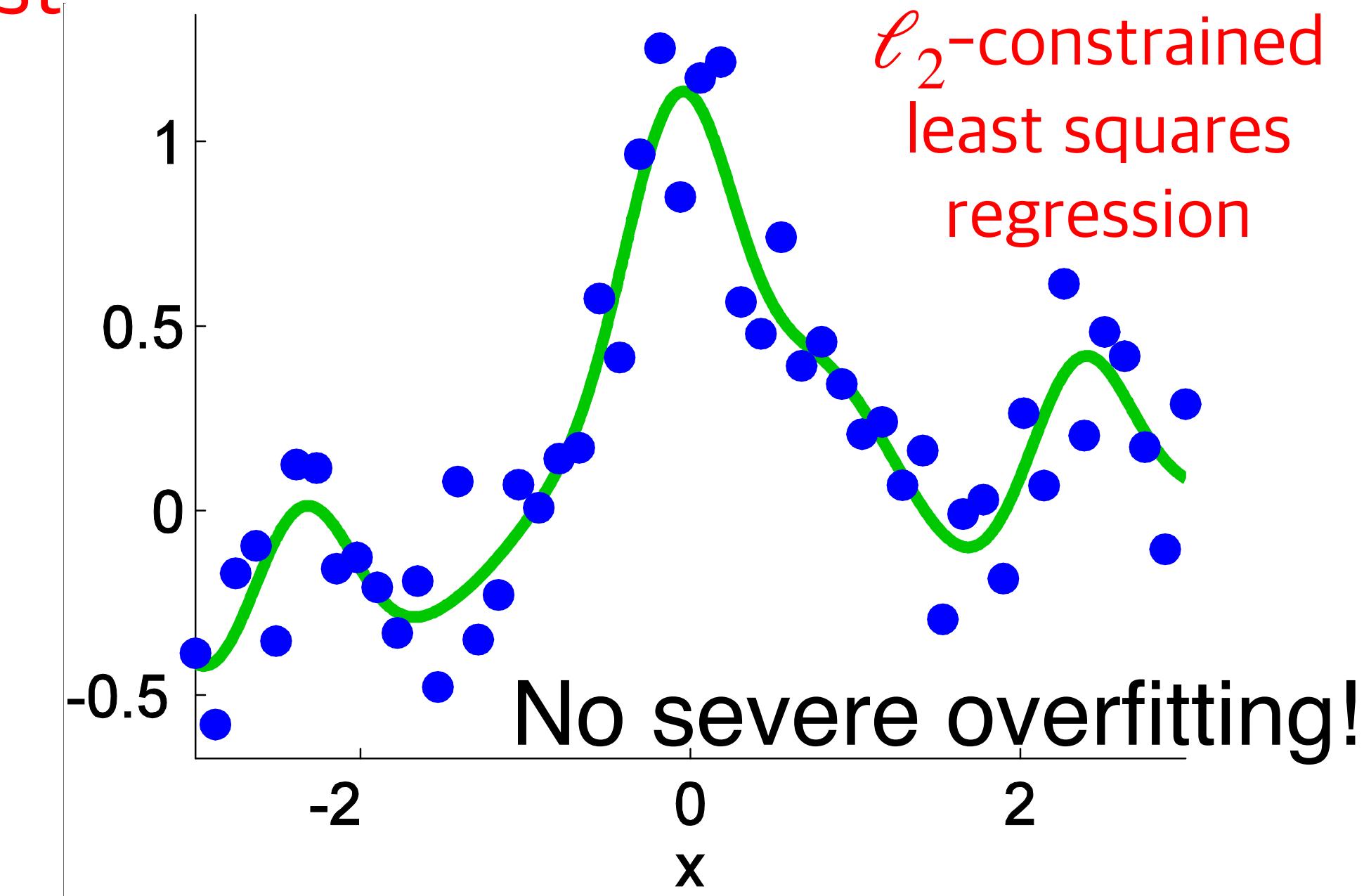
$$\hat{\theta} = (K^2 + \lambda I)^{-1} K^\top y$$



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

$$K = \begin{pmatrix} K(x_1, x_1) & \cdots & K(x_1, x_n) \\ \vdots & \ddots & \vdots \\ K(x_n, x_1) & \cdots & K(x_n, x_n) \end{pmatrix}$$

$$y = (y_1, \dots, y_n)^\top$$



# Example with Python

```
import numpy as np; import matplotlib; matplotlib.use('TkAgg')
import matplotlib.pyplot as plt; np.random.seed(0) # set seed for reproducibility

def generate_sample(xmin, xmax, sample_size):
    x = np.linspace(start=xmin, stop=xmax, num=sample_size)
    pix = np.pi * x
    target = np.sin(pix) / pix + 0.1 * x
    noise = 0.05 * np.random.normal(loc=0., scale=1., size=sample_size)
    return x, target + noise

def calc_design_matrix(x, c, h):
    return np.exp(-(x[None] - c[:, None]) ** 2 / (2 * h ** 2))

# create sample
sample_size = 50
xmin, xmax = -3, 3
x, y = generate_sample(xmin=xmin, xmax=xmax, sample_size=sample_size)
# calculate design matrix
h = 0.1
k = calc_design_matrix(x, x, h)
# solve the least square problem
l = 0.3
theta = np.linalg.solve(
    k.T.dot(k) + l * np.identity(len(k)),
    k.T.dot(y[:, None]))
# create data to visualize the prediction
X = np.linspace(start=xmin, stop=xmax, num=5000)
K = calc_design_matrix(x, X, h)
prediction = K.dot(theta)
# visualization
plt.clf(); plt.scatter(x, y, c='green', marker='o')
plt.plot(X, prediction); plt.show()
```

# Contents

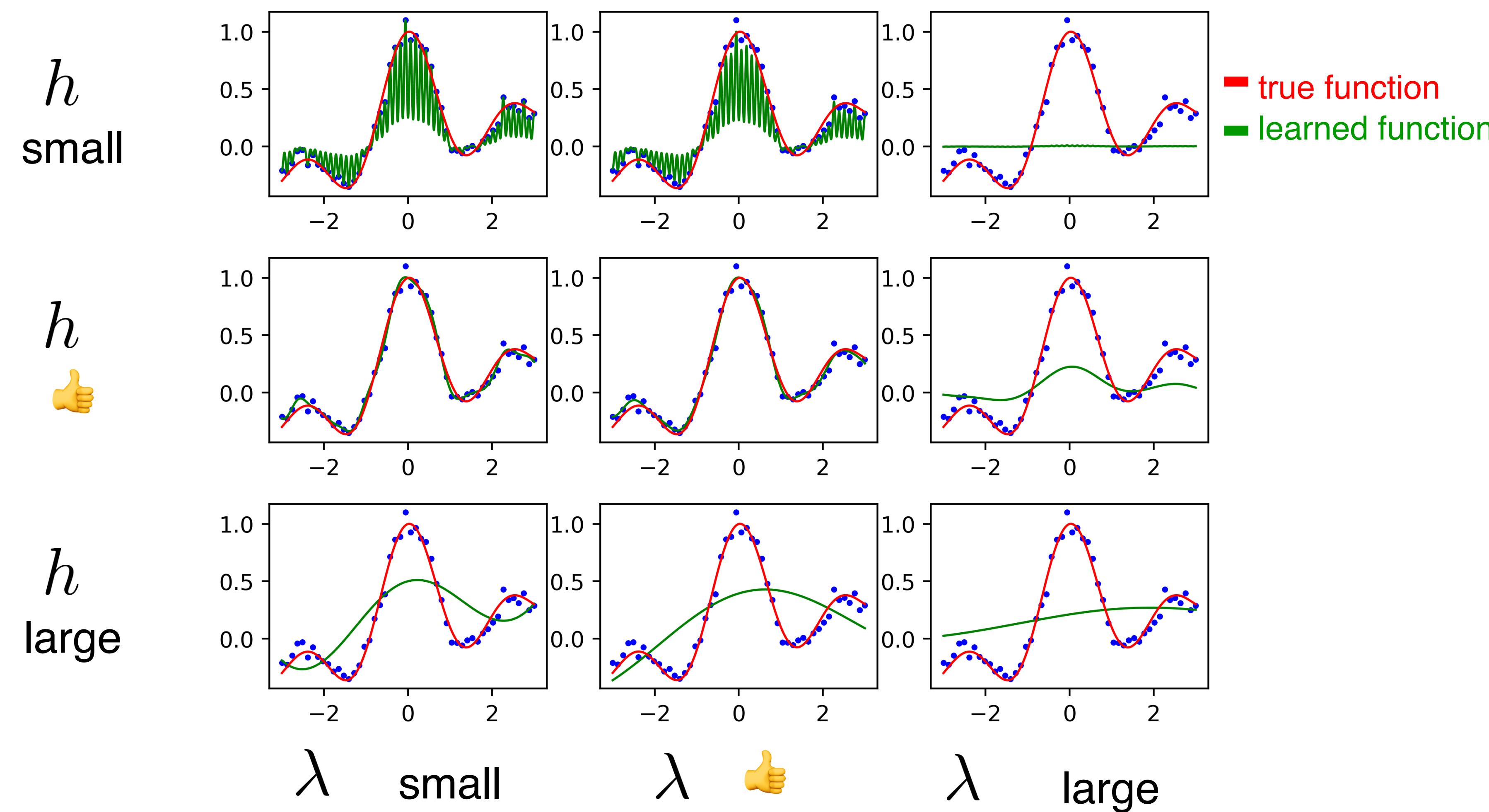
45

1. Learning model (Chap.2)
2. Least squares regression (Chap.3)
3. Regularized regression (Chap. 4)
  - A) Constrained least squares regression
  - B) Model selection

# Model selection w/ regularized regression

46

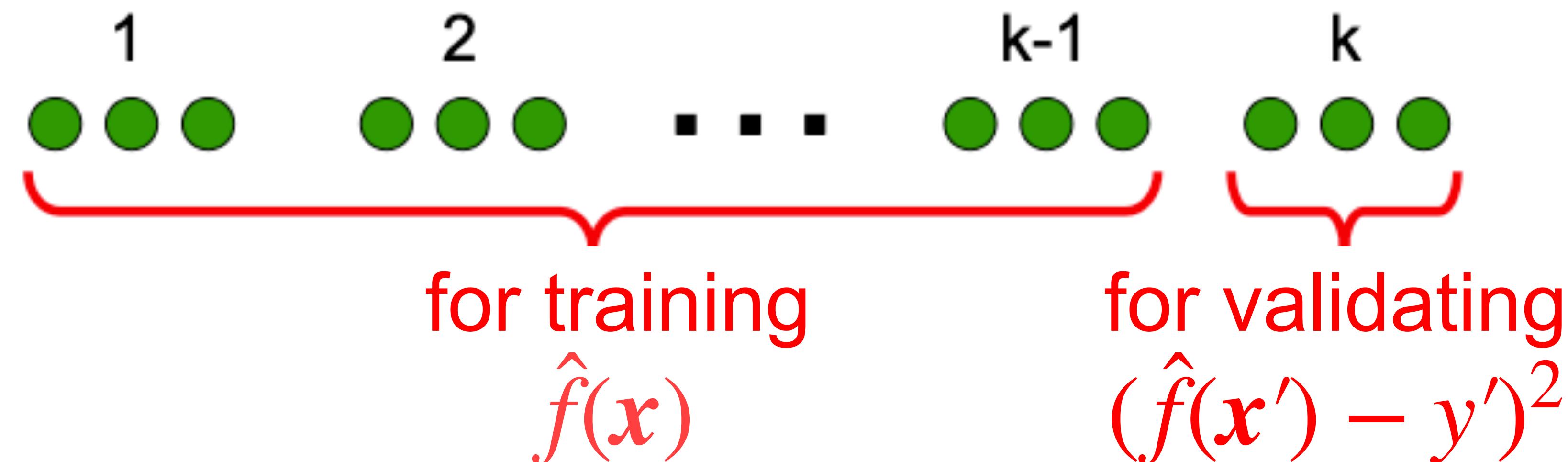
- The results of regularized regression depends on how you choose the regularization param  $\lambda$  and Gaussian bandwidth  $h$ .



# Cross validation

47

- Split training samples  $Z = \{(x_i, y_i)\}_{i=1}^n$  into  $k$  groups:  $\{Z_i\}_{i=1}^k$
- Use samples from groups excluding  $Z_i$  and learn  $\theta$  (fix  $\lambda, h$ ).
- Use the remaining  $Z_i$  to check the test error.
- Repeat this for all  $i \in [k]$ , and return the mean of the test errors.

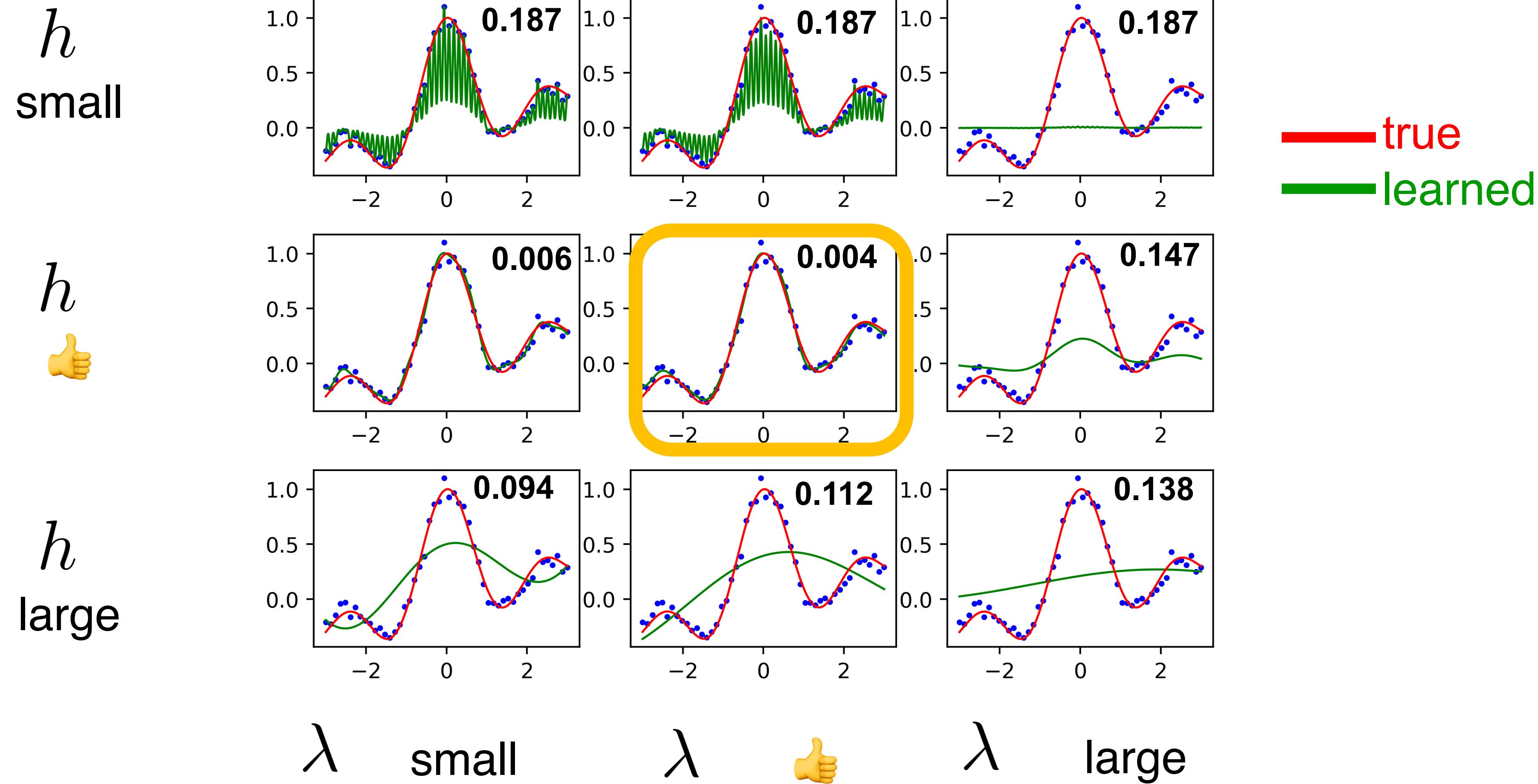


# Example of CV

48

- Gaussian kernel model:

$$f_{\theta}(x) = \sum_{j=1}^n \theta_j \exp\left(-\frac{\|x - x_j\|^2}{2h^2}\right)$$



Homework for this week is  
implementing this.

# Leave-one-out CV

49

- # of groups  $k = \# \text{ of samples } n$

$$\frac{1}{n} \sum_{i=1}^n \left( \hat{f}_i(\mathbf{x}_i) - y_i \right)^2$$

$\hat{f}_i(\mathbf{x})$ : learned from data other than  $(\mathbf{x}_i, y_i)$

- Leave one sample out each time
- In general, regularized regression can be computationally expensive and may not be practical.
- However, for linear models, there are some regularized regression techniques that can be computed analytically.

Proving this is homework.

$$\frac{1}{n} \|\tilde{\mathbf{H}}^{-1} \mathbf{H} \mathbf{y}\|^2$$

$$\mathbf{H} = \mathbf{I} - \Phi (\Phi^\top \Phi + \lambda \mathbf{I})^{-1} \Phi^\top$$

$\tilde{\mathbf{H}}$ : diagonals are same as  $\mathbf{H}$ , non-diagonals are zero

# Summary of regularized regression

50

- Least squares regression can easily overfit to noise.
- Constrain the search range of parameters using the  $\ell_2$ -norm.
- Solution can be analytically derived.
- Model selection is important.

# Contents

51

1. Learning model (Chap.2)
2. Least squares regression (Chap.3)
3. Regularized regression (Chap. 4)

# Summary

- Models for learning a function:
  - Linear models, kernel models, non-linear models
- Least squares regression:
  - Minimizes the squared error between training samples
  - The solution can be computed analytically
- Regularized regression:
  - Reduces overfitting in least squares regression
  - The solution can be computed analytically
  - Model selection is performed using cross-validation

# Schedule

- |    |       |                  |     |       |                              |
|----|-------|------------------|-----|-------|------------------------------|
| 1. | 04/8  | Introduction     | 8.  | 06/17 | Deep learning 3              |
| 2. | 04/15 | Regression 1     | 9.  | 06/24 | Semi-supervised<br>learning  |
| 3. | 04/22 | Regression 2     | 10. | 07/01 | Language models              |
| ●  | 04/30 | Cancelled        | 11. | 07/08 | Representation<br>learning 1 |
| 4. | 05/13 | Classification 1 | 12. | 07/15 | Representation<br>learning 2 |
| 5. | 05/20 | Classification 2 | 13. | 07/22 | Advanced topics              |
| 6. | 05/27 | Deep learning 1  |     |       |                              |
| ●  | 06/03 | No lecture       |     |       |                              |
| 7. | 06/10 | Deep learning 2  |     |       |                              |

# Homework 1

54

- Implement cross-validation for  $\ell_2$ -regularized least squares regression with a Gaussian kernel model and determine the regularization parameter and Gaussian width.

$$f_{\theta}(\mathbf{x}) = \sum_{j=1}^n \theta_j K(\mathbf{x}, \mathbf{x}_j)$$

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

# Homework 2

55

- Show that the mean squared error for leave-one-out cross-validation of  $\ell_2$ -regularized regression using a linear model  $f_{\theta}(\mathbf{x}) = \sum_{j=1}^b \theta_j \phi_j(\mathbf{x})$  can be analytically calculated as follows:

$$\frac{1}{n} \|\tilde{\mathbf{H}}^{-1} \mathbf{H} \mathbf{y}\|^2$$

$$\mathbf{H} = \mathbf{I} - \Phi (\Phi^\top \Phi + \lambda \mathbf{I})^{-1} \Phi^\top$$

$\tilde{\mathbf{H}}$ : diagonals are same as  $\mathbf{H}$ , non-diagonals are zero

# Homework 2 (Continued)

56

## ■ Helpful information:

- Express  $\hat{\theta}_i$  which is learned by removing  $(\mathbf{x}_i, y_i)$  with:

$$\phi_i = (\phi_1(\mathbf{x}_i), \phi_2(\mathbf{x}_i), \dots, \phi_b(\mathbf{x}_i))^\top$$

$$\Phi = (\phi_1, \phi_2, \dots, \phi_n)^\top,$$

$$\mathbf{y}, y_i, U = \Phi^\top \Phi + \lambda I$$

- A special case of Sherman–Morrison–Woodbury:

$$(U - uu^\top)^{-1} = U^{-1} + \frac{U^{-1}uu^\top U^{-1}}{1 - u^\top U^{-1}u}$$

# Homework 2 (Continued)

57

## ■ Helpful information:

- Express prediction for  $y_i, \phi_i^\top \hat{\theta}_i$  (where the parameter was learned by removing  $(x_i, y_i)$ ) with the following.

$$\phi_i = (\phi_1(x_i), \phi_2(x_i), \dots, \phi_b(x_i))^\top$$

$$\Phi = (\phi_1, \phi_2, \dots, \phi_n)^\top,$$

$$y, y_i, U = \Phi^\top \Phi + \lambda I$$

- Expanding on the denominator that appears in the inverse matrix formula and repeating the process of common denominator