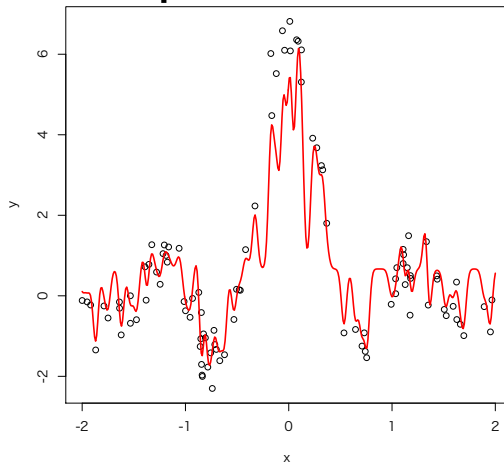


# 1 Regression Analysis

Training data:  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ .

$$y = \underbrace{f(\mathbf{x})}_{\text{function}} + \underbrace{\varepsilon}_{\text{r.v.: error}} \longrightarrow \text{estimate } f(\mathbf{x})$$

- Simple model is not good to learn complex data structure.  
→ Complex model is desirable
- Too complex models  $\implies$  overfitting



Overfit to data  
→ low prediction accuracy

- It is crucial to tuning the model complexity properly.

- Map  $\mathbf{x}$  to high-dimensional space:

$$\mathbf{x} \longmapsto \boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_D(\mathbf{x}))^T$$

$\phi_k(\mathbf{x})$ : (non-linear) basis functions

- Linear Regression Model:

$$f(\mathbf{x}) = \sum_{k=1}^D a_k \phi_k(\mathbf{x}) = \mathbf{a}^T \boldsymbol{\phi}(\mathbf{x})$$

Estimate the coefficient  $\mathbf{a}$  from training data.

- Choose functions  $\phi(\mathbf{x})$  having a “nice” property  
→ the computation is tractable.
- To avoid overfitting, regularization and cross validation are useful.

## — Kernel Regression Analysis —

- least square method with kernel-based modeling

# Estimation for Linear Regression Models

- Least Square Method (LSM):

$$\sum_{i=1}^n (y_i - \phi(\mathbf{x}_i)^T \mathbf{a})^2 = \|\mathbf{y} - \Phi^T \mathbf{a}\|^2 \rightarrow \text{minimize w.r.t. } \mathbf{a},$$

where  $\Phi = (\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)) \in \mathbb{R}^{D \times n}$ ,  $\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$ .

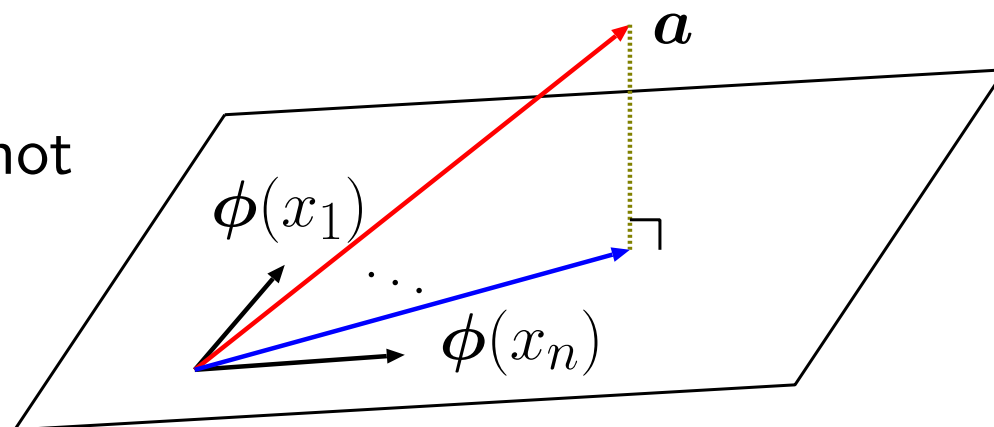
- $\text{rank } \Phi = D \Rightarrow \hat{\mathbf{a}} = (\Phi \Phi^T)^{-1} \Phi \mathbf{y}$ .

Another expression of the solution:

$$\min_{\mathbf{a} \in \mathbb{R}^D} \sum_{i=1}^n (y_i - \phi(\mathbf{x}_i)^T \mathbf{a})^2$$

- the solution lies on  $\text{span}\{\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)\}$ .

Orthogonal component does not affect the square error.



- $\mathbf{a} = \sum_{j=1}^n \beta_j \phi(\mathbf{x}_j) = \Phi \boldsymbol{\beta}, \quad \boldsymbol{\beta} \in \mathbb{R}^n.$

$$\sum_{i=1}^n (y_i - \phi(\mathbf{x}_i)^T \mathbf{a})^2 = \|\mathbf{y} - \Phi^T \Phi \boldsymbol{\beta}\|^2 \longrightarrow \min_{\boldsymbol{\beta}}$$

Optimality conditions :  $\Phi^T \Phi \underbrace{\Phi \hat{\boldsymbol{\beta}}}_{\hat{\mathbf{a}}} = \Phi^T \Phi \mathbf{y}$

Define  $n$  by  $n$  matrix  $K = (K_{ij})$  as

$$K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) \stackrel{\text{def}}{=} \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \in \mathbb{R},$$

$$\implies K = \Phi^T \Phi$$

- $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$  is called kernel function
- $K$ : Gram matrix

(the rigorous definition is given later)

- Optimality condition:

$$\begin{aligned}\Phi^T \Phi \Phi^T \Phi \hat{\beta} &= \Phi^T \Phi y \iff K^2 \hat{\beta} = Ky \\ \implies \text{calculate } \hat{\beta} &= (\hat{\beta}_1, \dots, \hat{\beta}_n)^T\end{aligned}$$

- Estimated regression function:  $\hat{f}(x) = \phi(x)^T \hat{a}$ .

$$\hat{f}(x) = \phi(x)^T \underbrace{\sum_{i=1}^n \phi(x_i) \hat{\beta}_i}_{\hat{a}} = \sum_{i=1}^n k(x, x_i) \hat{\beta}_i$$

- kernel function  $k(x, x') \implies$  estimator  $\hat{f}(x)$

Examples of kernel functions:  $\mathbf{x} \in \mathbb{R}^d \mapsto \phi(\mathbf{x}) \in \mathbb{R}^D$ .

- linear kernel:  $D = d$ .

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}', \quad (\phi(\mathbf{x}) = \mathbf{x})$$

Model:  $y = \mathbf{a}^T \phi(\mathbf{x}) + \varepsilon = \mathbf{a}^T \mathbf{x} + \varepsilon$



- Polynomial kernel of degree  $\ell \in \mathbb{N}$ :  $D = \frac{(\ell+d)!}{\ell! d!}$

$$k(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^\ell,$$

Model:  $y = \mathbf{a}^T \boldsymbol{\phi}(\mathbf{x}) + \varepsilon$ .

$\boldsymbol{\phi}(\mathbf{x})$ : all monomials of degree  $\leq \ell$ .

For  $d = 2$ ,  $\ell = 2$  and  $\mathbf{x} = (x_1, x_2)^T \in \mathbb{R}^2$ ,

$$\boldsymbol{\phi}(\mathbf{x}) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2)^T.$$

$$\begin{aligned} \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{z}) &= 1 + 2x_1z_1 + 2x_2z_2 + x_1^2z_1^2 + x_2^2z_2^2 + 2x_1x_2z_1z_2 \\ &= (1 + x_1z_1 + x_2z_2)^2 \end{aligned}$$

- Gaussian kernel:  $D = \infty$ .

$$k(\mathbf{x}, \mathbf{x}') = \exp\{-\sigma \cdot \|\mathbf{x} - \mathbf{x}'\|^2\}, \quad \sigma > 0$$

For  $d = 1$ ,  $\sigma = 1$  and  $x \in \mathbb{R}$ ,

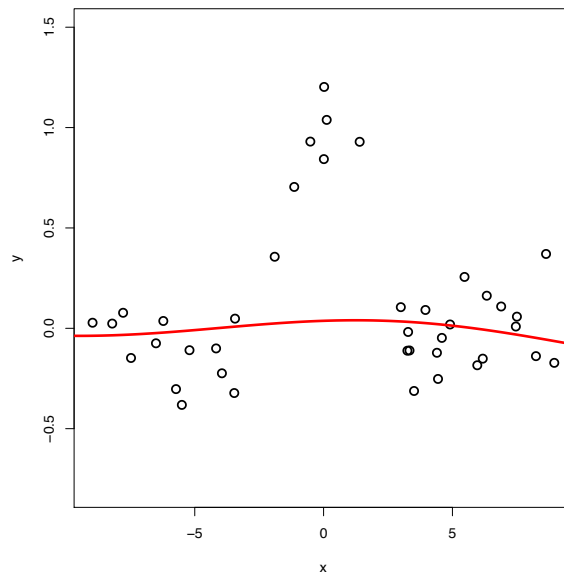
$$\boldsymbol{\phi}(x) = (\phi_0(x), \phi_1(x), \phi_2(x), \dots)^T, \quad \phi_j(x) = \frac{x^j e^{-x^2/2}}{\sqrt{j!}}, \quad x \in \mathbb{R}$$

# Overfitting to data

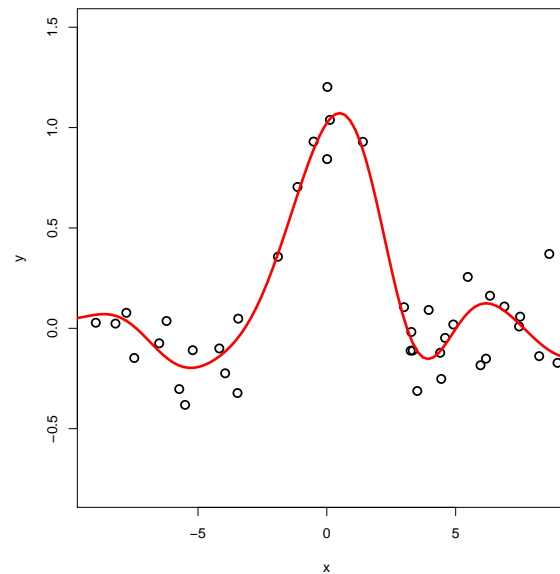
- simple model: hard to deal with complex data  
→ use the model with many parameters
- model with too many parameters does not work.  
overfitting to data.

degree of freedom

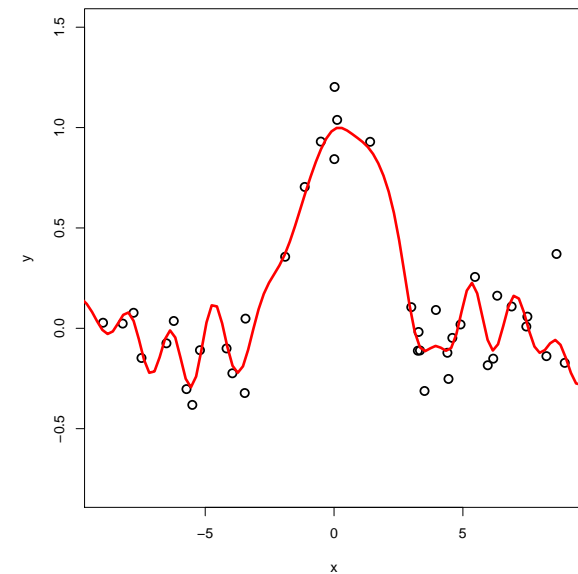
small



medium



large



# Regularization: tune the degree of freedom

large model & appropriate constraint

linear regression model :  $y = \underline{\mathbf{a}^T \phi(\mathbf{x}) + b} + \varepsilon$

ex.  $\phi(x) = (x, x^2, x^3, \dots, x^{100})$ ,  $\phi(\mathbf{x})$  of Gaussian kernel, etc.

data:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ .

$$\min_{\mathbf{a}, b} \sum_{i=1}^n (y_i - (\phi(\mathbf{x}_i)^T \mathbf{a} + b))^2 + \underbrace{\frac{\lambda \|\mathbf{a}\|^2}{2}}_{\text{regularization term}} \quad (\text{Ridge regression})$$

$$\implies \text{opt. sol. } \hat{\mathbf{a}}, \hat{b}. \quad \hat{f}(x) = \hat{\mathbf{a}}^T \phi(\mathbf{x}) + \hat{b}$$

$$\min_{\mathbf{a} \in \mathbb{R}^D, b \in \mathbb{R}} \|\mathbf{y} - \Phi^T \mathbf{a} - b \mathbf{1}\|^2 + \lambda \|\mathbf{a}\|^2, \quad \Phi = (\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n))$$

regularization parameter  $\lambda > 0$ .

$\lambda$ : large

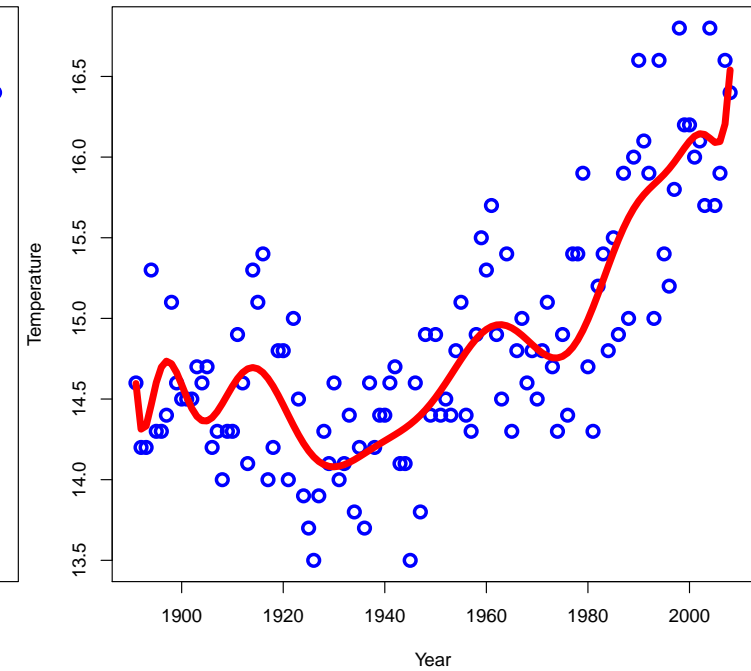
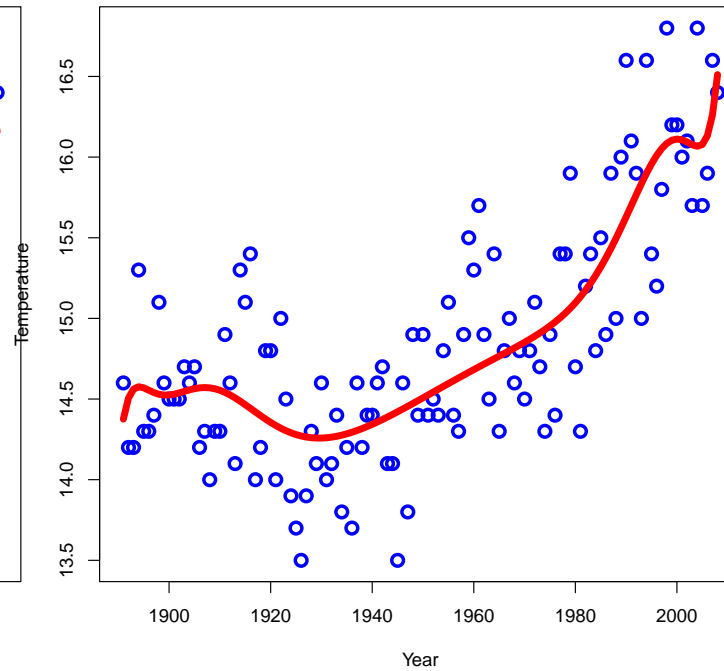
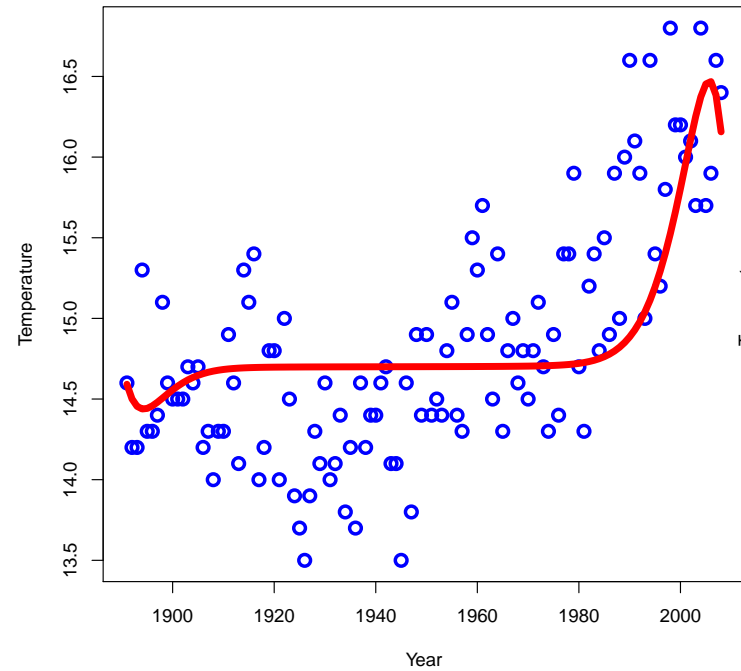
$\lambda$ : medium

$\lambda$ : small

lambda=10000

lambda=1

lambda=1e-05



small  $\longleftarrow$  degree of freedom  $\longrightarrow$  large

Kernel representation of Ridge regression:

$$\min_{\mathbf{a}, b} \|\mathbf{y} - \Phi^T \mathbf{a} - b \mathbf{1}\|^2 + \lambda \|\mathbf{a}\|^2, \quad \Phi = (\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)).$$

In the same way as the standard LMS, the optimal  $\hat{\mathbf{a}}$  lies on the subspace  $\text{span}\{\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)\}$ .

Substitute  $\mathbf{a} = \sum_{i=1}^n \phi(\mathbf{x}_i) \beta_i = \Phi \boldsymbol{\beta}$ , then for  $K = \Phi^T \Phi$ ,

$$\begin{aligned} \|\mathbf{y} - \Phi^T \mathbf{a} - b \mathbf{1}\|^2 + \lambda \|\mathbf{a}\|^2 &= \|\mathbf{y} - \Phi^T \Phi \boldsymbol{\beta} - b \mathbf{1}\|^2 + \lambda \boldsymbol{\beta}^T \Phi^T \Phi \boldsymbol{\beta} \\ &= \|\mathbf{y} - K \boldsymbol{\beta} - b \mathbf{1}\|^2 + \lambda \boldsymbol{\beta}^T K \boldsymbol{\beta} \longrightarrow \min_{\boldsymbol{\beta}, b} \end{aligned}$$

# Kernel-Ridge Regression

- Optimality condition

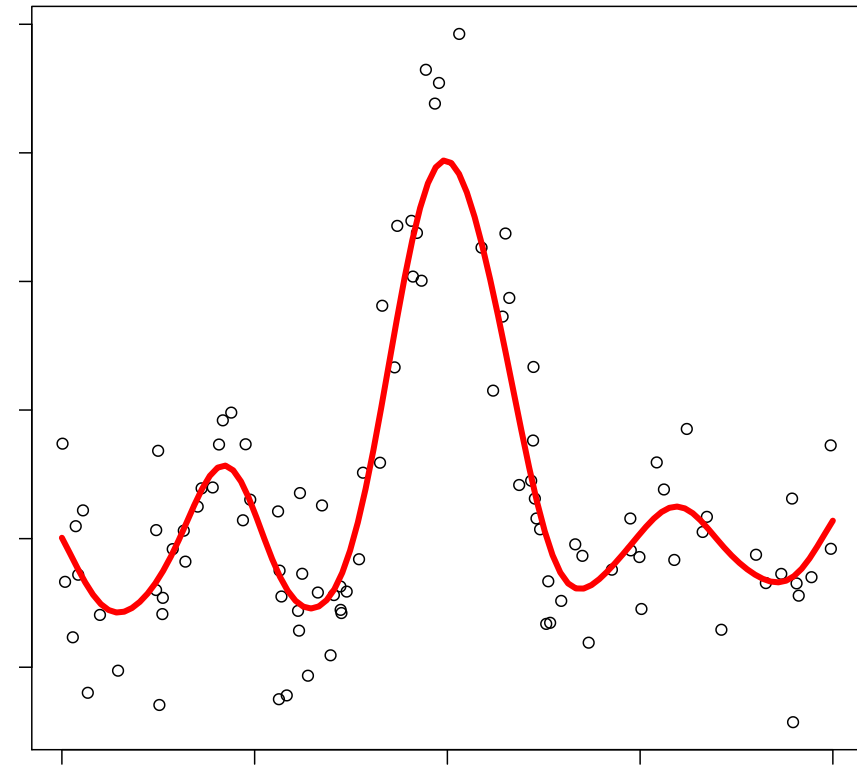
$$\min_{\boldsymbol{\beta}, b} \|\mathbf{y} - K\boldsymbol{\beta} - b\mathbf{1}\|^2 + \lambda\boldsymbol{\beta}^T K\boldsymbol{\beta}$$
$$\implies \begin{pmatrix} K + \lambda I & \mathbf{1} \\ \mathbf{1}^T K & n \end{pmatrix} \begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{b} \end{pmatrix} = \begin{pmatrix} \mathbf{y} \\ \mathbf{1}^T \mathbf{y} \end{pmatrix}$$

- estimated regression function:

$$\hat{f}(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^T \sum_{i=1}^n \boldsymbol{\phi}(\mathbf{x}_i) \hat{\beta}_i + \hat{b} = \sum_{i=1}^n k(\mathbf{x}, \mathbf{x}_i) \hat{\beta}_i + \hat{b}$$

Plot: estimated regression function

- kernel width:  $\sigma = 3$
- regularization par.:  $\lambda = 1$





## — Model Selection —

How to choose regularization parameter  $\lambda$  ?

- Training error and Test error
- Cross Validation for model parameter tuning

# Kernel-Ridge Regression

Gaussian kernel :  $k(\mathbf{x}, \mathbf{x}') = \exp\{-\sigma \cdot \|\mathbf{x} - \mathbf{x}'\|^2\}$

We need to determine the following **model parameters**:

- Regularization par.:  $\lambda$
- kernel parameter:  $\sigma$

How to choose  $\lambda$  and  $\sigma$  ?

Note. For the polynomial kernel, we need to determine  $\lambda$  and the degree  $\ell$ .

# Training error and Test error

- training data:  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \sim_{i.i.d.} p(\mathbf{x}, y)$
- estimated regression function  $\hat{f}(\mathbf{x})$

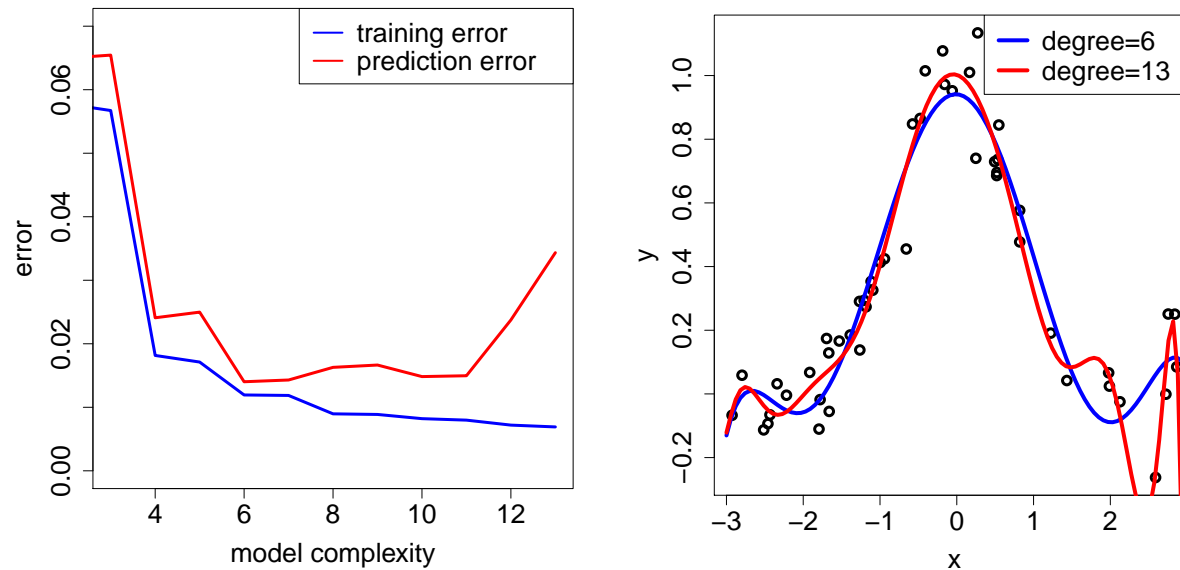
training error of  $\hat{f}(\mathbf{x})$  :  $\frac{1}{n} \sum_{i=1}^n (\hat{f}(\mathbf{x}_i) - y_i)^2$  (calculated from data)

test error of  $\hat{f}(\mathbf{x})$  :  $\mathbb{E}_{(\mathbf{x}, y) \sim P} [(\hat{f}(\mathbf{x}) - y)^2]$  ( $P$  is unknown)

———— Purpose of Regression Analysis ————

Find  $\hat{f}(\mathbf{x})$  that achieves a small (or minimum) test error.

## Polynomial regression



### overfitting

- For high degree polynomial models, we have
  - small training error
  - large test error
  - overfitting: large gap between training error and test error.

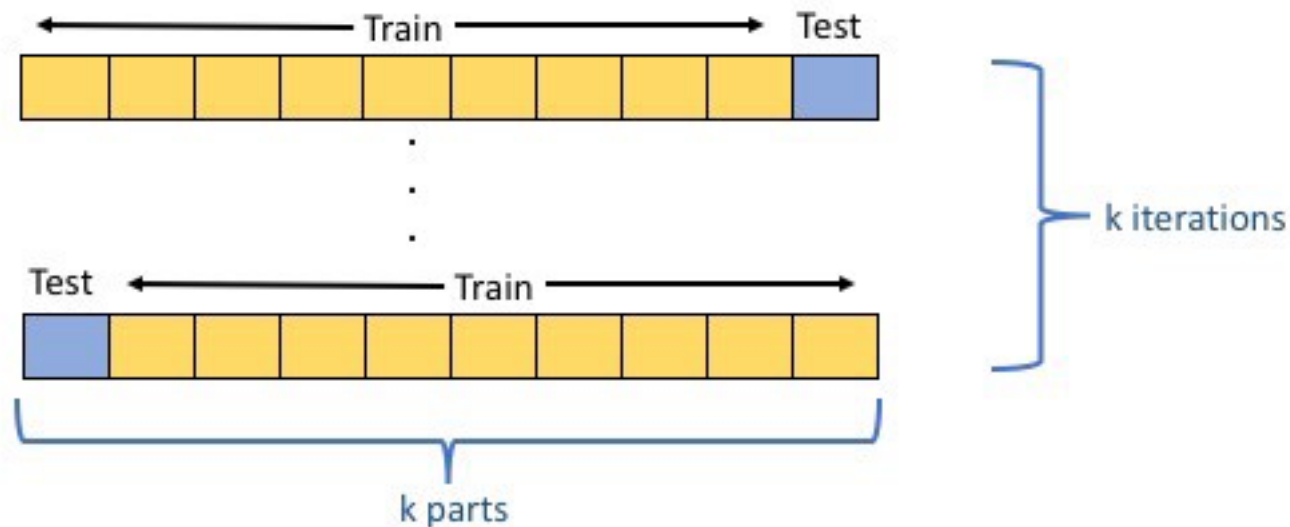
appropriate model complexity is required.

**Cross validation:** estimator of test error

## *K*-fold Cross Validation Method

Fixing a model parameter, say  $\lambda$  and  $\sigma$ , execute the following procedure.

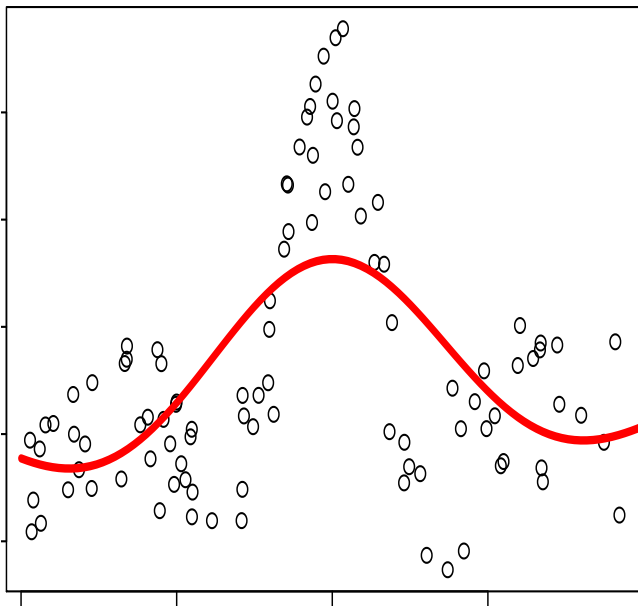
1. Divide the training data into  $k$  parts.
2. Use  $k - 1$  of the parts for training, and 1 for testing.
3. Repeat the procedure  $k$  times, rotating the test set.
4. Calculate an expected performance metric (mean square error/test error rate) based on the results across the iterations



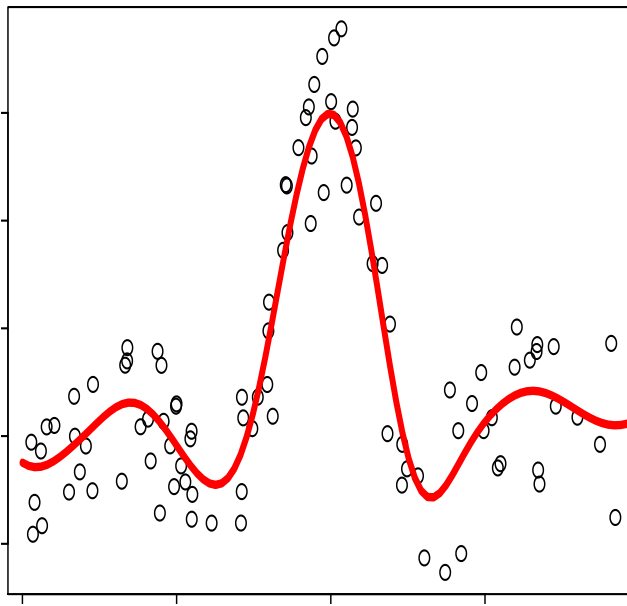
<https://medium.com/@mtterribile/understanding-cross-validations-purpose-53490faf6a86>

# Example: Kernel regression with Gaussian kernel

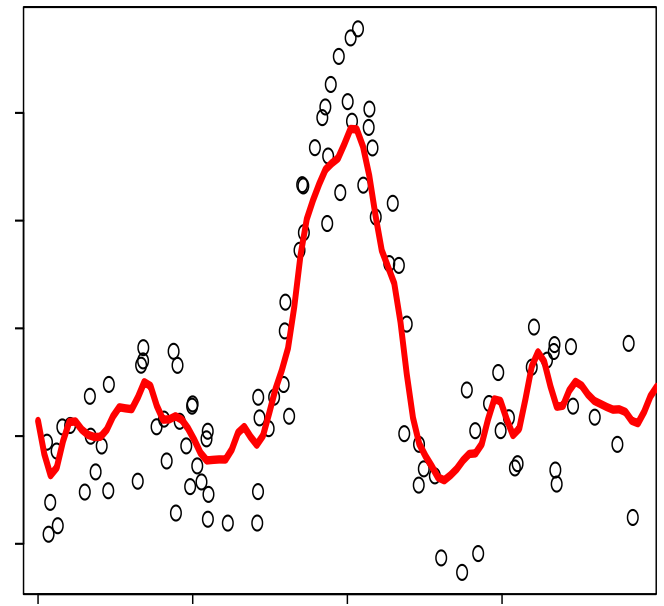
- kernel par.  $\sigma > 0$  is determined by  $K$ -cv.
- regularization par.  $\lambda > 0$  is fixed to 1.



$$(\sigma, \lambda) = (0.1, 1)$$



$$(\sigma, \lambda) = (2.5, 1)$$



$$(\sigma, \lambda) = (30, 1)$$

# Median Heuristics for Gaussian Kernel

Gaussian kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp\{-\sigma \|\mathbf{x} - \mathbf{x}'\|^2\}$$

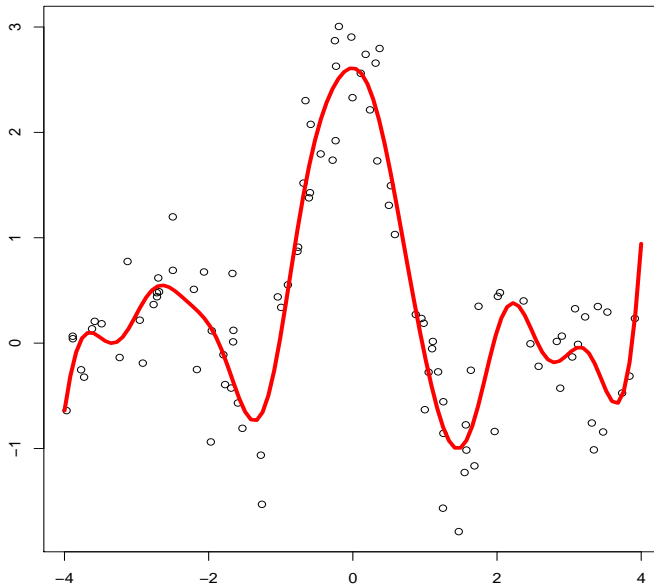
- For computational stability, choose  $\sigma$  such that

$\sigma \|\mathbf{x}_i - \mathbf{x}_j\|^2$  takes values around 1.

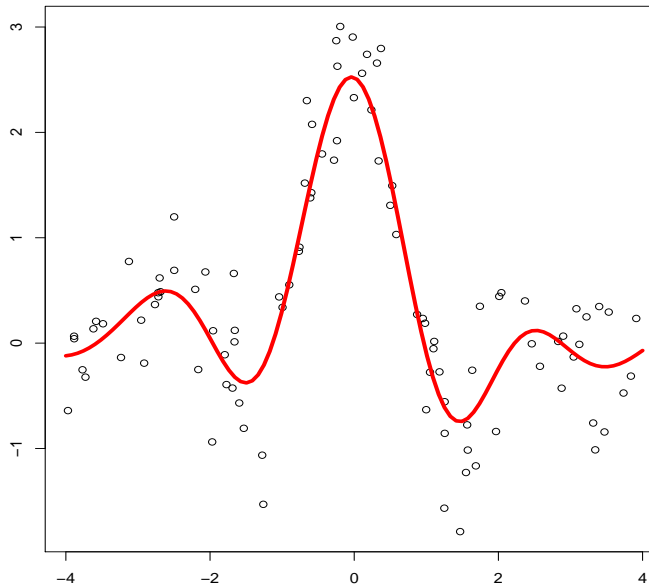
$$\sigma \longleftarrow \frac{1}{\text{median}\{\|\mathbf{x}_i - \mathbf{x}_j\|^2 \mid i < j\}}$$

# Example

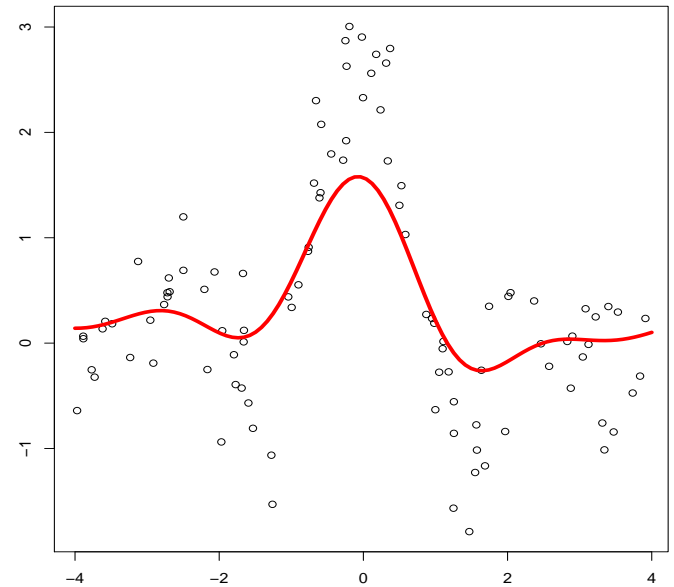
- kernel parameter  $\sigma$  is determined by the heuristics.
- regularization par.  $\lambda$ :  $K$ -cv



$$(\text{sigma}, \lambda) = (0.98, 0.1^5)$$



$$(\text{sigma}, \lambda) = (0.98, 0.78)$$



$$(\text{sigma}, \lambda) = (0.98, 10)$$