

# Regression Analysis I

**Jin Young Choi**

**Seoul National University**

# Outline

---

- linear regression
  - simple linear regression
  - multiple linear regression
- nonlinear regression
  - logistic regression
  - high-order regression
  - basis-function regression
- matrix form for regression
  - recursive least squares
- partial least squares
  - over-fitting and underfitting
  - bias/variance
  - principle component regression
  - partial least squares algorithm
  - ridge regression
  - lasso, elastic regression
- Gaussian process regression

# LINEAR REGRESSION

---

**JIN YOUNG CHOI**

**ECE, SEOUL NATIONAL UNIVERSITY**

<http://3.droppdf.com/files/pjxkl/regression-analysis-by-example-5th-edition.pdf>

<https://github.com/jwangjie/Gaussian-Processes-Regression-Tutorial>

# Regression Analysis

---

- For independent random variable  $X$ , and dependent random variable  $Y$ , assume they have a functional correlation between them, i.e.

$$Y = f(X)$$

- **Regression**: a process to find a parametric model  $\hat{f}$  that gives the best fit of  $f$  for the observed samples

$$Y = \hat{f}(X) + \epsilon, \quad X: \text{predictor r.v.}, Y: \text{response r.v.}$$

- Assume  $E(\epsilon) = 0$ ,  $\text{var}(\epsilon) = \sigma^2$ , then  $E(Y|x) = \hat{f}(x)$  for an observed non-random value  $x$
- $\hat{f}$  can be estimated from the sample pairs  $\{(y_i, x_i) | i = 1, 2, \dots, n\}$

$$y_i = \hat{f}(x_i) + \epsilon_i, \quad i = 1, \dots, n,$$

where  $\epsilon_i$  are i.i.d. zero mean and variance  $\sigma^2$

# Simple Linear Regression

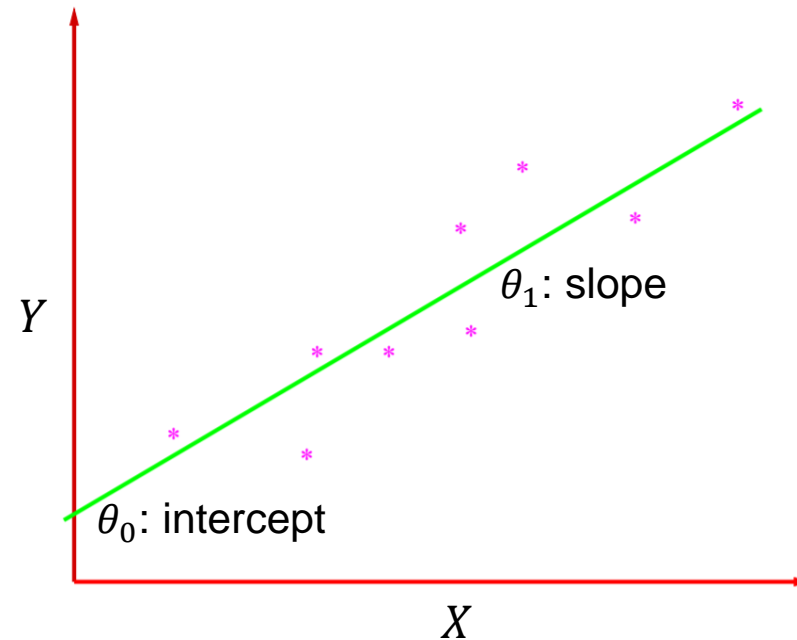
- Simple linear regression model

$$Y = \theta_0 + \theta_1 X + \epsilon$$

$$y_i = \theta_0 + \theta_1 x_i + \epsilon_i, \quad i = 1, \dots, n,$$

where  $\theta_0$ : intercept,  $\theta_1$ : slope

Observation Number	Response Y	Predictor X
1	$y_1$	$x_1$
2	$y_2$	$x_2$
3	$y_3$	$x_3$
$\vdots$	$\vdots$	$\vdots$
$n$	$y_n$	$x_n$



# Simple Linear Regression

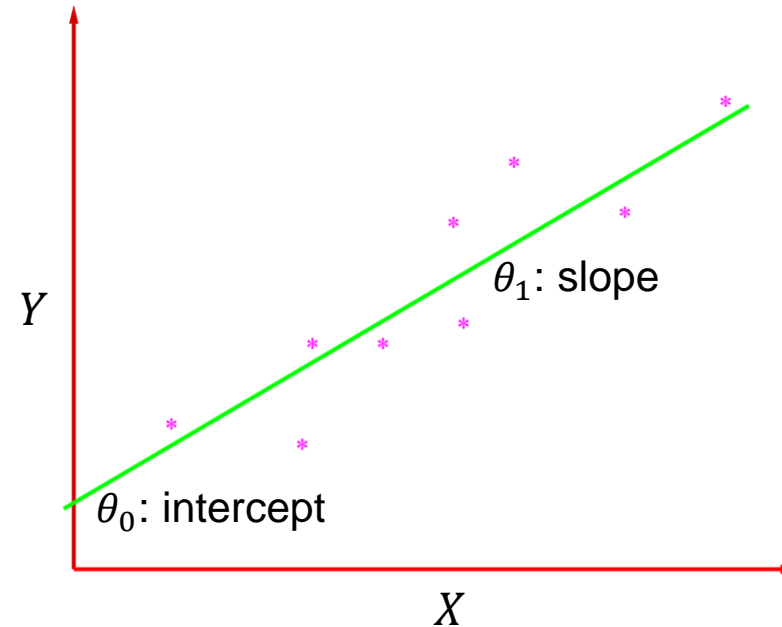
- Correlation of  $Y$  &  $X$

$$Y = \theta_0 + \theta_1 X + \epsilon$$

$$\text{Cov}(Y, X) = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y}) (x_i - \bar{x})$$

where  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ ,  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Observation Number	Response Y	Predictor X
1	$y_1$	$x_1$
2	$y_2$	$x_2$
3	$y_3$	$x_3$
$\vdots$	$\vdots$	$\vdots$
n	$y_n$	$x_n$



# Simple Linear Regression

- Correlation of  $Y$  &  $X$

$$Y = \theta_0 + \theta_1 X + \epsilon$$

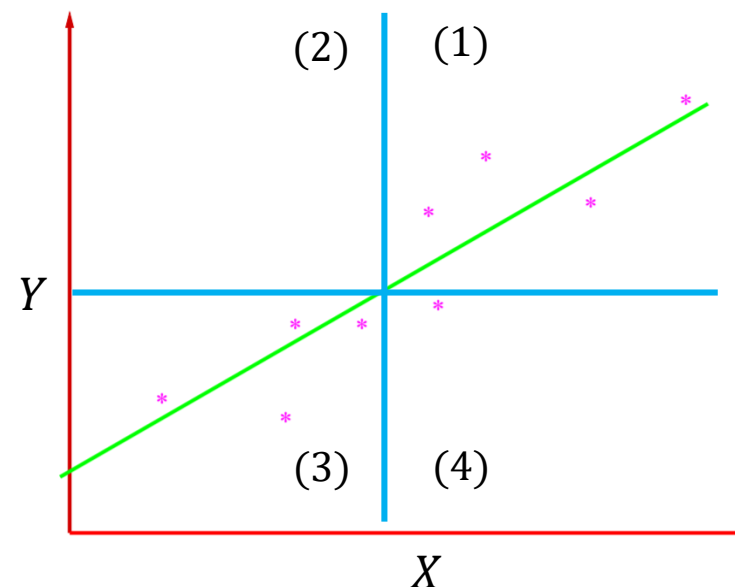
$$\text{Cov}(Y, X) = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})$$

where  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ ,  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Q	$y_i - \bar{y}$	$x_i - \bar{x}$	$(y_i - \bar{y})(x_i - \bar{x})$
(1)	+	+	+
(2)	+	-	-
(3)	-	-	+
(4)	-	+	-

$$\theta_1 \geq 0 \quad \longrightarrow \quad \text{Cov}(Y, X) \geq 0$$

$$\theta_1 < 0 \quad \longrightarrow \quad \text{Cov}(Y, X) < 0$$



# Simple Linear Regression

- Correlation Coefficient of  $Y$  &  $X$

$$Y = \theta_0 + \theta_1 X + \epsilon$$

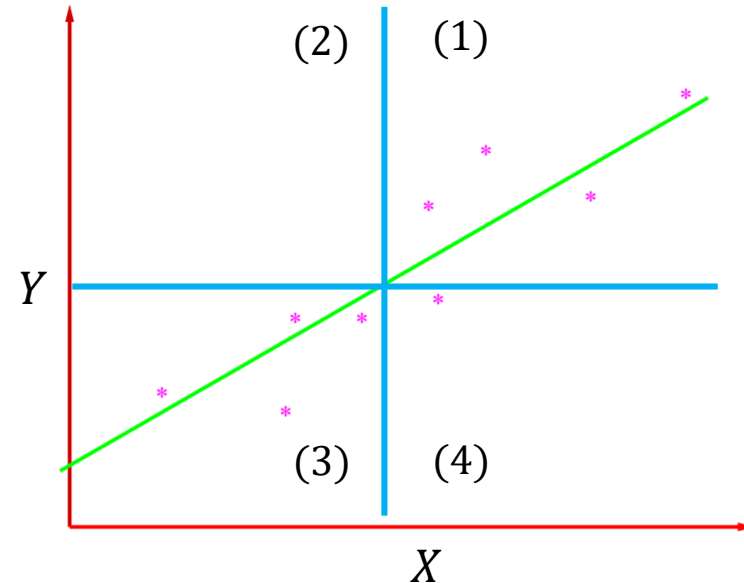
$$\rho(Y, X) = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \bar{y}}{\sigma_y} \right) \left( \frac{x_i - \bar{x}}{\sigma_x} \right)$$

where  $\sigma_y^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$ ,  $\sigma_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$

Q	$y_i - \bar{y}$	$x_i - \bar{x}$	$(y_i - \bar{y})(x_i - \bar{x})$
(1)	+	+	+
(2)	+	-	-
(3)	-	-	+
(4)	-	+	-

$$\theta_1 \geq 0 \quad \longrightarrow \quad 1 \geq \rho(Y, X) \geq 0$$

$$\theta_1 < 0 \quad \longrightarrow \quad -1 \leq \rho(Y, X) < 0$$





# Parameter Estimation

---

- Least Squares Estimation

Parameters are estimated by maximum likelihood estimation (MLE)

$$\epsilon_i = y_i - \theta_0 + \theta_1 x_i, \quad i = 1, \dots, n, \quad \epsilon_i \sim N(0, \sigma^2)$$

MLE:

$$(\hat{\theta}_0, \hat{\theta}_1) = \operatorname{argmax}_{(\theta_0, \theta_1)} \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\epsilon_i^2}{2\sigma^2}\right)$$

$$(\hat{\theta}_0, \hat{\theta}_1) = \operatorname{argmax}_{(\theta_0, \theta_1)} \ln \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\epsilon_i^2}{2\sigma^2}\right)$$

$$(\hat{\theta}_0, \hat{\theta}_1) = \operatorname{argmin}_{(\theta_0, \theta_1)} \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

LSE:

$$\text{minimizing} \quad S(\theta_0, \theta_1) = \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2.$$

Solution:

$$\text{by } \partial S / \partial \theta_0 = 0, \partial S / \partial \theta_1 = 0 \text{ at } \hat{\theta}_0 \text{ \& \& } \hat{\theta}_1,$$

# Maximum Likelihood Estimation

$$\theta^* = \operatorname{argmax}_{\theta} p(\{\epsilon_i\}|\theta)$$

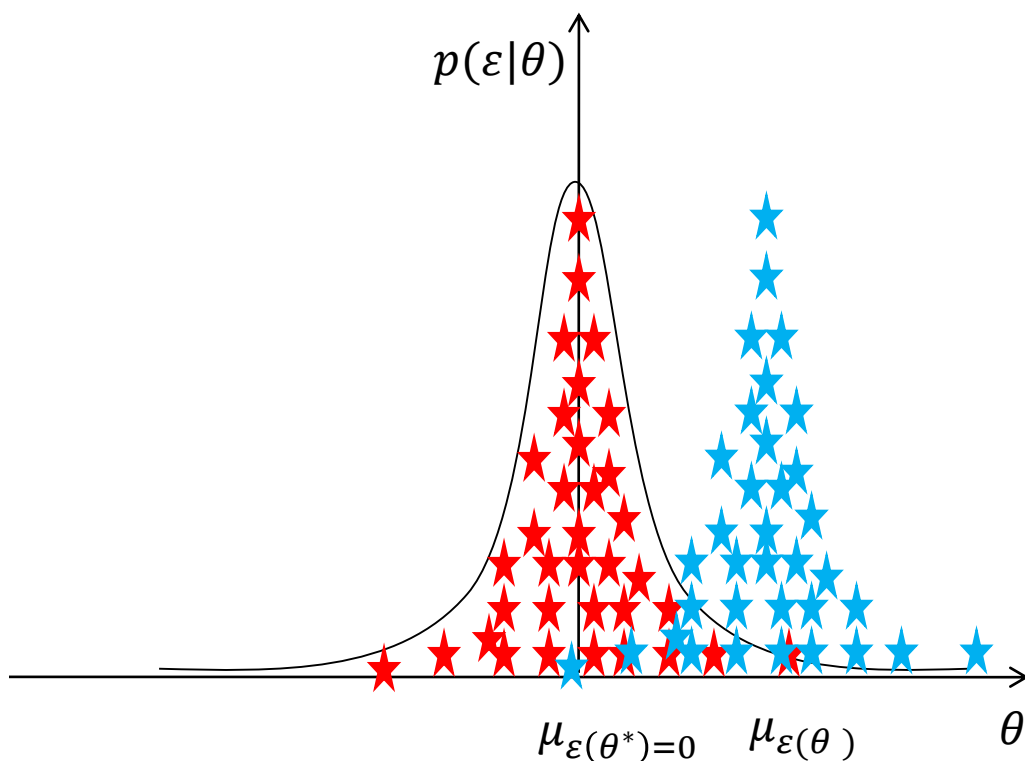
$$\epsilon(\theta) = Y - \theta X$$

$$(\hat{\theta}_0, \hat{\theta}_1) = \operatorname{argmax}_{(\theta_0, \theta_1)} \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\epsilon_i^2}{2\sigma^2}\right)$$

$$(\hat{\theta}_0, \hat{\theta}_1) = \operatorname{argmax}_{(\theta_0, \theta_1)} \ln \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\epsilon_i^2}{2\sigma^2}\right)$$

$$(\hat{\theta}_0, \hat{\theta}_1) = \operatorname{argmin}_{(\theta_0, \theta_1)} \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

$$\hat{\theta} = \operatorname{argmin}_{\theta} \|\epsilon\|^2 = \|\mathbf{y} - \Phi\theta\|^2 \cong S(\theta)$$



$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \phi_1^T \\ \phi_2^T \\ \vdots \\ \phi_n^T \end{bmatrix} \theta + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}, \quad \Phi_k = \begin{bmatrix} \phi_{11} & \phi_{12} & \cdots & \phi_{1p} \\ \phi_{21} & \phi_{22} & \cdots & \phi_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{n1} & \phi_{n2} & \cdots & \phi_{np} \end{bmatrix}$$

$$y_i = \phi_i^T \theta + \epsilon_i$$

$$y_i = \theta_0 + \theta_1 \phi_{i1} + \theta_2 \phi_{i2} + \cdots + \theta_p \phi_{i(p-1)} + \epsilon_i,$$

$$i = 1, \dots, n$$

# Parameter Estimation

---

- Least Squares Estimation

Parameters are estimated by maximum likelihood estimation (MLE)

$$\epsilon_i = y_i - \theta_0 + \theta_1 x_i, \quad i = 1, \dots, n, \quad \epsilon_i \sim N(0, \sigma^2)$$

MLE:

$$(\hat{\theta}_0, \hat{\theta}_1) = \operatorname{argmax}_{(\theta_0, \theta_1)} \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\epsilon_i^2}{2\sigma^2}\right)$$

$$(\hat{\theta}_0, \hat{\theta}_1) = \operatorname{argmax}_{(\theta_0, \theta_1)} \ln \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\epsilon_i^2}{2\sigma^2}\right)$$

$$(\hat{\theta}_0, \hat{\theta}_1) = \operatorname{argmin}_{(\theta_0, \theta_1)} \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

LSE:

$$\text{minimizing} \quad S(\theta_0, \theta_1) = \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2.$$

Solution:

$$\text{by } \partial S / \partial \theta_0 = 0, \partial S / \partial \theta_1 = 0 \text{ at } \hat{\theta}_0 \text{ \& \& } \hat{\theta}_1,$$

# Parameter Estimation

- Least Squares Estimation

$$\epsilon_i = y_i - \theta_0 + \theta_1 x_i, \quad i = 1, \quad \dots, \quad n.$$

LSE:

$$(\hat{\theta}_0, \hat{\theta}_1) = \underset{(\theta_0, \theta_1)}{\operatorname{argmin}} S(\theta_0, \theta_1) = \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2.$$

Solution:

by  $\partial S / \partial \theta_0 = 0, \partial S / \partial \theta_1 = 0$  at  $\hat{\theta}_0$  &  $\hat{\theta}_1$ ,

$$\sum_{i=1}^n (y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i) = 0, \quad \rightarrow \quad \boxed{\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}}$$

$$\sum_{i=1}^n (y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i) x_i = 0, \rightarrow \sum_{i=1}^n (y_i - \bar{y} - \hat{\theta}_1 (x_i - \bar{x})) (x_i - \bar{x} + \bar{x}) = 0,$$

$$\rightarrow \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) - \hat{\theta}_1 \sum_{i=1}^n (x_i - \bar{x})^2 = 0 \rightarrow \boxed{\hat{\theta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

# Parameter Estimation

- Least squares regression line

$$\hat{Y} = \hat{\theta}_0 + \hat{\theta}_1 X.$$

Fitted values:

$$\hat{y}_i = \hat{\theta}_0 + \hat{\theta}_1 x_i, \quad i = 1, \dots, n.$$

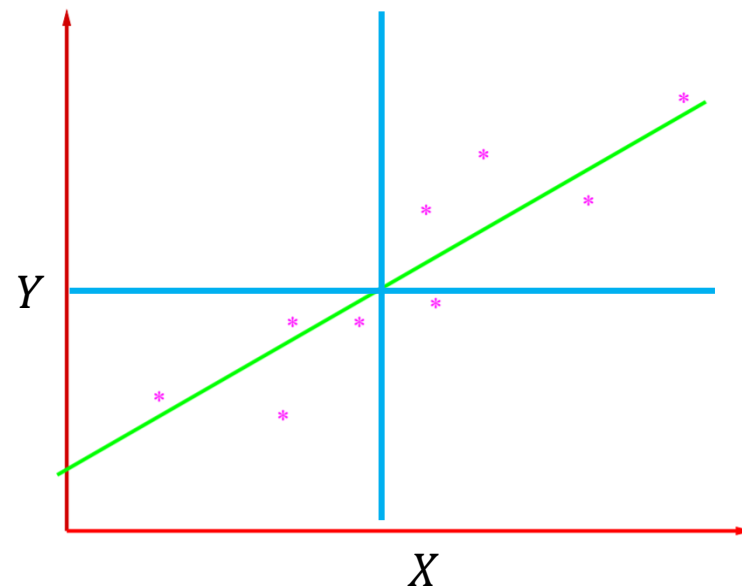
Error to the  $i$ -th observation:

$$e_i = y_i - \hat{y}_i, \quad i = 1, \dots, n.$$

Alternative formula for  $\hat{\theta}_1$ :

$$\hat{\theta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\text{Cov}(Y, X)}{\text{Var}(X)} = \frac{\rho(Y, X)\sigma_x\sigma_y}{\sigma_x^2} = \rho(Y, X) \frac{\sigma_y}{\sigma_x}$$

→ slope has the **same sign** with the correlation ( $\rho(Y, X)$ ; covariance)



# Measuring the Quality of Fit

- Original Model:

$$Y = \theta_0 + \theta_1 X + \epsilon.$$

Least squares regression line:

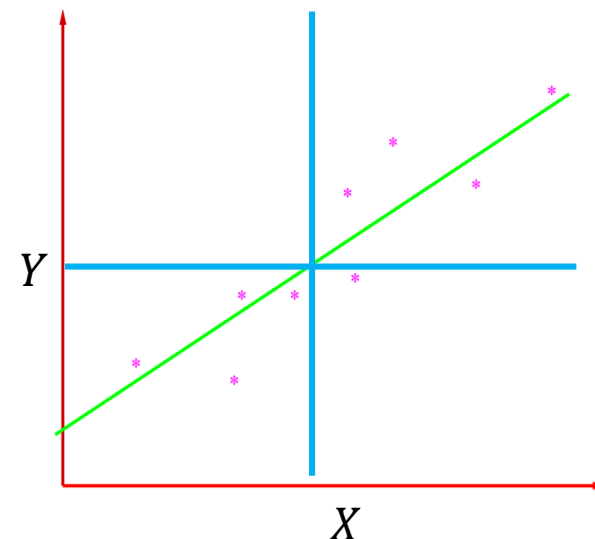
$$\hat{Y} = \hat{\theta}_0 + \hat{\theta}_1 X.$$

- Correlation between  $Y$  &  $\hat{Y}$  :

$$\rho(Y, \hat{Y}) = \frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\left(\sum_{i=1}^n (y_i - \bar{y})^2 \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2\right)}}$$

Note that  $\rho(Y, \hat{Y})$  can not be negative. Why?

Note that  $\rho(Y, \hat{Y}) = 1$  implies the perfect fit.



# Measuring the Quality of Fit

- Goodness-of-fit index:

$SST: \sum_{i=1}^n (y_i - \bar{y})^2$ ,  $SST$ : Total sum of squares

$SSR: \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ ,  $SSR$ : Regression (explained) sum of squares

$SSE: \sum_{i=1}^n (y_i - \hat{y}_i)^2$ ,  $SSE$ : Residual (error) sum of squares

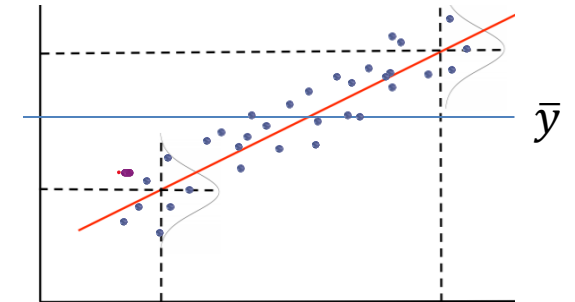
- Interpretation:

$$\begin{array}{rclcl}
 y_i & = & \hat{y}_i & + & y_i - \hat{y}_i \\
 \text{Observed} & = & \text{Fit} & + & \text{Error} \\
 y_i - \bar{y} & = & \hat{y}_i - \bar{y} & + & y_i - \hat{y}_i \\
 \text{Deviation} & & \text{Deviation to Fit} & & \text{Residual}
 \end{array}$$

$$SST = SSR + SSE \quad \because \sum_{i=1}^n (\hat{y}_i - \bar{y})(y_i - \hat{y}_i) = 0 \quad [1]$$

- $R^2$ : Coefficient of determination

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} \quad (R = 1 \text{ implies the perfect fit})$$



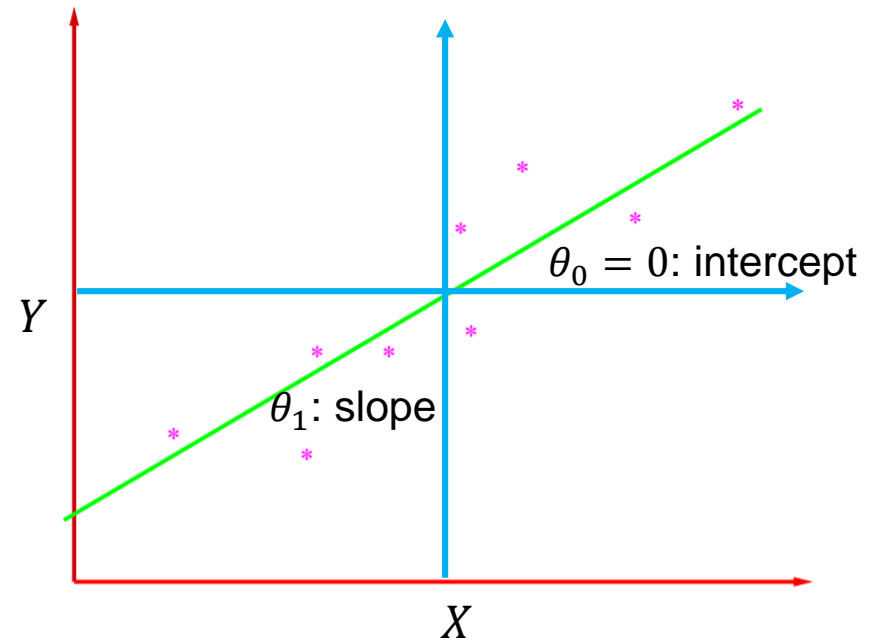
# Regression Line through Origin

- Simple linear regression model

$$Y = \theta_0 + \theta_1 X + \epsilon$$

$$Y = \theta_1 X + \epsilon, \quad \text{no-intercept model, } \bar{y} = \bar{x} = 0$$

Observation Number	Response Y	Predictor X
1	$y_1 - \bar{y}$	$x_1 - \bar{x}$
2	$y_2 - \bar{y}$	$x_2 - \bar{x}$
3	$y_3 - \bar{y}$	$x_3 - \bar{x}$
$\vdots$	$\vdots$	$\vdots$
n	$y_n - \bar{y}$	$x_n - \bar{x}$





# Regression Line through Origin

---

- *no-intercept* model

$$y_i = \theta_1 x_i + \epsilon,$$

$$\hat{y}_i = \hat{\theta}_1 x_i, \quad i = 1, \dots, n$$

$$e_i = y_i - \hat{y}_i.$$

$$\text{Cov}(Y, X) = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) \rightarrow \text{Cov}(Y, X) = \frac{1}{n} \sum_{i=1}^n y_i x_i$$

$$\rho(Y, X) = \frac{1}{n} \sum_{i=1}^n \frac{y_i x_i}{\sigma_y \sigma_x}, \quad \sigma_y^2 = \frac{1}{n} \sum_{i=1}^n y_i^2, \quad \sigma_x^2 = \frac{1}{n} \sum_{i=1}^n x_i^2$$

$$\hat{\theta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \rightarrow \hat{\theta}_1 = \frac{\sum_{i=1}^n y_i x_i}{\sum_{i=1}^n x_i^2} = \frac{\text{Cov}(Y, X)}{\sigma_x^2} = \rho(Y, X) \frac{\sigma_y}{\sigma_x}$$

$$R^2 = \frac{\sum_{i=1}^n \hat{y}_i^2}{\sum_{i=1}^n y_i^2} = 1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n y_i^2}$$

# Multivariate Linear Regression

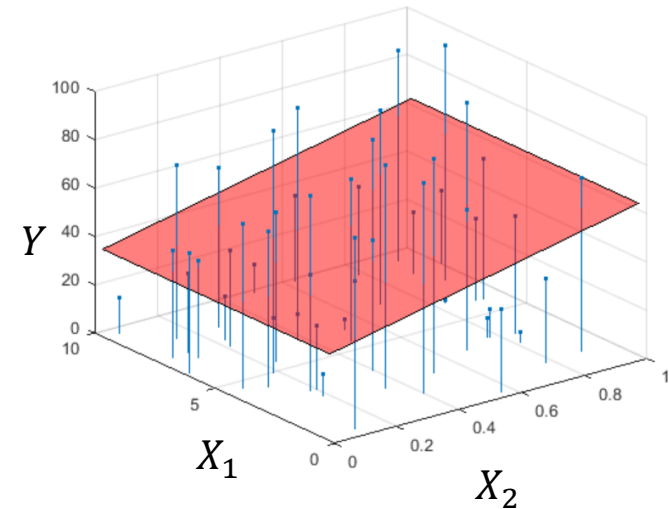
- **Multivariate** linear regression model:  $p$  predictor (explanatory) variables

$$Y = \theta_0 + \theta_1 X_1 + \theta_2 X_2 + \cdots + \theta_p X_p + \epsilon$$

$$y_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2} + \cdots + \theta_p x_{ip} + \epsilon_i, \quad i = 1, \dots, n,$$

where  $\theta_0$ : **intercept**,  $(\theta_1, \theta_2, \dots, \theta_p)$ : normal vector (ex.;  $y = w^T x + b$ )

$i$	$Y$	Predictor			
		$X_1$	$X_2$	$\cdots$	$X_p$
1	$y_1$	$x_{11}$	$x_{12}$	$\cdots$	$x_{1p}$
2	$y_2$	$x_{21}$	$x_{22}$	$\cdots$	$x_{2p}$
3	$y_3$	$x_{31}$	$x_{32}$	$\cdots$	$x_{3p}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$y_n$	$x_{n1}$	$x_{n2}$	$\cdots$	$x_{np}$



# Multivariate Linear Regression

- Multivariate linear regression model:  $p$  predictor (explanatory) variables

$$Y = \theta_0 + \theta_1 X_1 + \theta_2 X_2 + \cdots + \theta_p X_p + \epsilon$$

$$y_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2} + \cdots + \theta_p x_{ip} + \epsilon_i, \quad i = 1, \dots, n,$$

where  $\theta_0$ : intercept,  $(\theta_1, \theta_2, \dots, \theta_p)$ : normal vector

- Fitted model by LSE:  $n - p - 1$ ; degree of freedom ( $df$ );  $p + 1$ ; # of estimated parameters

$$\hat{y}_i = \hat{\theta}_0 + \hat{\theta}_1 x_{i1} + \hat{\theta}_2 x_{i2} + \cdots + \hat{\theta}_p x_{ip}, \quad i = 1, \dots, n,$$

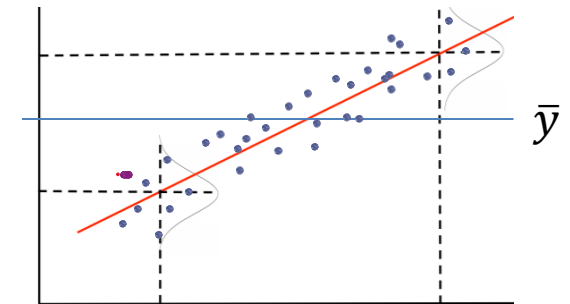
$$e_i = y_i - \hat{y}_i.$$

- Measuring Quality of Fit:

$$\rho(Y, \hat{Y}) = \frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{(\sum_{i=1}^n (y_i - \bar{y})^2 \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2)}}$$

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

- Adjusted  $R^2$ :  $R_a^2 = 1 - \frac{1/(n-p-1) \sum_{i=1}^n e_i^2}{1/(n-1) \sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{n-1}{n-p-1} (1 - R^2)$



# Multivariate Linear Regression

- Tests of Hypotheses for Multivariate linear model

$$Y = \theta_0 + \theta_1 X_1 + \theta_2 X_2 + \cdots + \theta_p X_p + \epsilon$$

$$y_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2} + \cdots + \theta_p x_{ip} + \epsilon_i, \quad i = 1, \dots, n,$$

where  $\theta_0$ : intercept,  $(\theta_1, \theta_2, \dots, \theta_p)$ : normal vector

- Hypotheses:  $H_0$ : Reduced model (RM),  $H_1$ : Full model (FM)

1. All the regression coefficients associated with the predictor variables are zero.
2. Some of the regression coefficients are zero.
3. Some of the regression coefficients are equal to each other.
4. The regression parameters satisfy certain specified constraints (ex.  $|\theta_i| \leq \alpha$ ).

- Sum of Squares:  $SSE(RM) \geq SSE(FM)$

$$SSE(FM) = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$SSE(RM) = \sum_{i=1}^n (y_i - \hat{y}_i^*)^2$$

- F-test:  $F = \frac{[SSE(RM) - SSE(FM)] / (p+1-k)}{SSE(FM) / (n-p-1)}$  ( $F$  is large  $\rightarrow$  RM is inadequate<sup>†</sup>)

<sup>†</sup> The critical values are given in Table A.4 and A.5 in "Regression Analysis by Example", S. Chatterjee et.al., Wiley.

# NONLINEAR REGRESSION

---

**JIN YOUNG CHOI**

**ECE, SEOUL NATIONAL UNIVERSITY**

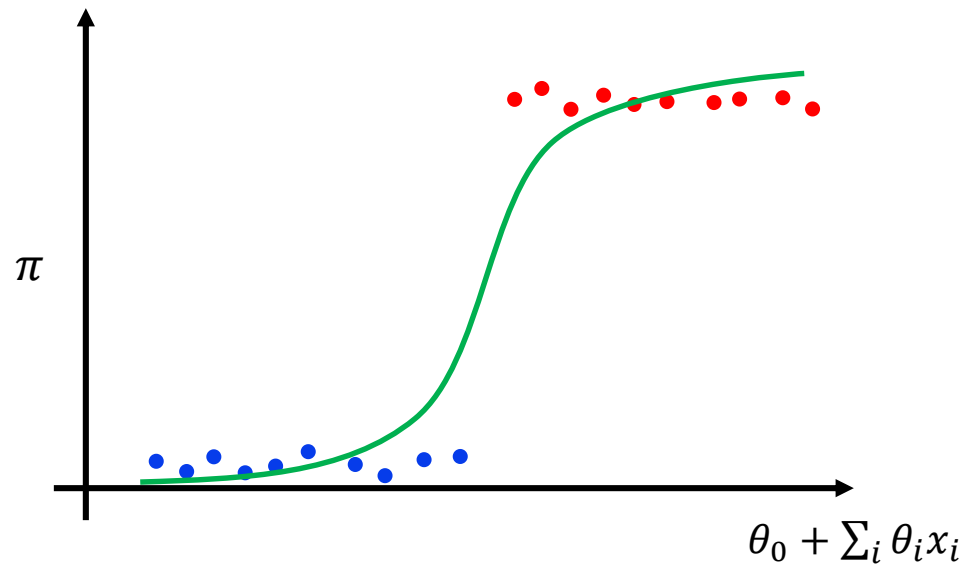
<https://github.com/jwangjie/Gaussian-Processes-Regression-Tutorial>

# Logistic Regression

---

- **Logistic response function** representing the relation between the probability  $\pi$  and  $X_1, X_2, \dots, X_p$

$$\pi = p(Y = 1 | X_1 = x_1, \dots, X_p = x_p) = \frac{\exp(\theta_0 + \theta_1 x_1 + \dots + \theta_p x_p)}{1 + \exp(\theta_0 + \theta_1 x_1 + \dots + \theta_p x_p)}$$



# Logistic Regression

- Logistic response function

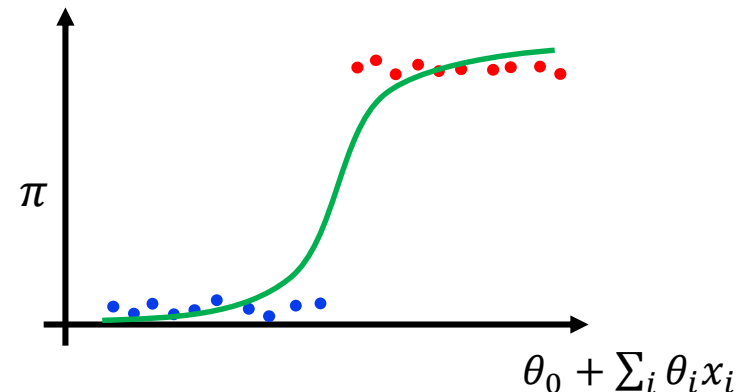
$$\pi(X_1 = x_1, \dots, X_p = x_p) = p(Y = 1 | X_1 = x_1, \dots, X_p = x_p) = \frac{\exp(\theta_0 + \theta_1 x_1 + \dots + \theta_p x_p)}{1 + \exp(\theta_0 + \theta_1 x_1 + \dots + \theta_p x_p)}$$

$$1 - \pi(X_1 = x_1, \dots, X_p = x_p) = p(Y = 0 | X_1 = x_1, \dots, X_p = x_p) = \frac{1}{1 + \exp(\theta_0 + \theta_1 x_1 + \dots + \theta_p x_p)}$$

$$\frac{\pi}{1 - \pi} = \exp(\theta_0 + \theta_1 x_1 + \dots + \theta_p x_p)$$

$$f(X_1 = x_1, \dots, X_p = x_p) = \ln \frac{\pi}{1 - \pi} = \theta_0 + \theta_1 x_1 + \dots + \theta_p x_p$$

$$f(X) = \ln \frac{\pi}{1 - \pi} = \theta_0 + \theta_1 X_1 + \dots + \theta_p X_p$$

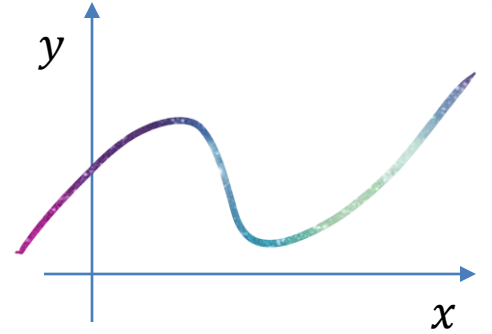


# High-order Regression

- High-order polynomial regression model

$$Y = \theta_0 + \theta_1 X + \theta_2 X^2 + \cdots + \theta_m X^m + \epsilon$$

$$y_i = \theta_0 + \theta_1 x_i + \theta_2 x_i^2 + \cdots + \theta_m x_i^m + \epsilon_i, \quad i = 1, \dots, n.$$



- High-order multivariate regression model

$$Y = \theta_0 + \theta_1 X_1 + \cdots + \theta_k X_k + \cdots + \theta_{k(\pi)} X_{\pi_1} \cdots X_{\pi_j} \cdots + \cdots + \theta_p X_{\mu(m)}^m + \epsilon$$

$$y_i = \theta_0 + \theta_1 x_{i1} + \cdots + \theta_k x_{ik} + \cdots + \theta_{k(\pi)} x_{i\pi_1} \cdots x_{i\pi_j} \cdots + \cdots + \theta_M x_{ip}^m + \epsilon_i$$

- Matrix-vector form

$$\text{Let } \theta = [\theta_0 \ \theta_1 \ \cdots \ \theta_p]^T, \quad \phi_i = [1 \ \phi_{i1} \ \cdots \ \phi_{ip}]^T$$

$$y = [y_1 \ y_2 \ \cdots \ y_n]^T, \quad \epsilon = [\epsilon_1 \ \epsilon_2 \ \cdots \ \epsilon_n]^T$$

$$\text{Then } y_i = \phi_i^T \theta + \epsilon_i, \quad i = 1, \dots, n.$$

$$y = \Phi \theta + \epsilon, \quad \Phi = [\phi_1 \ \phi_2 \ \cdots \ \phi_n]^T$$

$$\Phi = \begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^m \\ 1 & x_2 & x_2^2 & \cdots & x_2^m \\ 1 & x_3 & x_3^2 & \cdots & x_3^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^m \end{bmatrix}$$



# Basis-function Regression

---

- Matrix-vector form of General Regression

$$\text{Let } \theta = [\theta_0 \ \theta_1 \ \cdots \ \theta_p]^T, \ \phi_i = [1 \ \phi_{i1} \ \cdots \ \phi_{ip}]^T$$
$$\mathbf{y} = [y_1 \ y_2 \ \cdots \ y_n]^T, \ \boldsymbol{\epsilon} = [\epsilon_1 \ \epsilon_2 \ \cdots \ \epsilon_n]^T$$

$$\text{Then } y_i = \phi_i^T \theta + \epsilon_i, \ i = 1, \ \cdots, \ n.$$

$$\mathbf{y} = \Phi \theta + \boldsymbol{\epsilon}, \ \Phi = [\phi_1 \ \phi_2 \ \cdots \ \phi_n]^T$$

- Basis for General Regression

- sin, cos basis:  $\phi_{im} = \sin \omega_m x_i$  or  $\cos \omega_m x_i$

- radial basis:  $\phi_{im} = \exp \frac{-\|x_i - \mu_m\|^2}{\sigma_m^2}$

- sigmoid basis:  $\phi_{im} = \frac{1}{1 + \exp(-w_m^T x_i - b_m)}$  or  $\frac{\exp(w_m^T x_i + b_m)}{1 + \exp(w_m^T x_i + b_m)}$

Logistic Regression

# Parameter Estimation in Matrix form

---

- Least Squares Estimation

$$\mathbf{y} = \Phi\theta + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(0, \sigma^2 \mathbf{I})$$

MLE:

$$\hat{\theta} = \operatorname{argmax}_{\theta} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\|\boldsymbol{\epsilon}\|^2}{2\sigma^2}\right)$$

LSE:

$$\hat{\theta} = \operatorname{argmin}_{\theta} \|\boldsymbol{\epsilon}\|^2 = \|\mathbf{y} - \Phi\theta\|^2 \triangleq S(\theta)$$

Solution:

$$\text{by } \nabla_{\theta} S(\theta) = 0 \text{ at } \hat{\theta}.$$

# Parameter Estimation in Matrix form

---

- Least Squares Estimation

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \|\epsilon\|^2 = \|\mathbf{y} - \Phi\theta\|^2 \cong S(\theta)$$

Solution:

$$\nabla_{\theta} S(\theta) = 0 \text{ at } \hat{\theta}$$

$$\nabla_{\theta} (\mathbf{y} - \Phi\theta)^T (\mathbf{y} - \Phi\theta) = 0 \text{ at } \hat{\theta}$$

$$2\Phi^T (\mathbf{y} - \Phi\hat{\theta}) = 0$$

$$\Phi^T \mathbf{y} - \Phi^T \Phi \hat{\theta} = 0$$

$$\hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

# Interim Summary

---

- linear regression
  - simple linear regression
  - multiple linear regression
- nonlinear regression
  - logistic regression
  - high-order regression
  - basis-function regression
- matrix form for regression
  - recursive least squares
- partial least squares
  - over-fitting and underfitting
  - bias/variance
  - principle component regression
  - partial least squares algorithm
  - ridge regression
  - lasso, elastic regression
- Gaussian process regression

# Parameter Estimation in Matrix form

---

- Least Squares Estimation

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \|\epsilon\|^2 = \|\mathbf{y} - \Phi\theta\|^2 \cong S(\theta)$$

Solution:

$$\nabla_{\theta} S(\theta) = 0 \text{ at } \hat{\theta}$$

$$\nabla_{\theta} (\mathbf{y} - \Phi\theta)^T (\mathbf{y} - \Phi\theta) = 0 \text{ at } \hat{\theta}$$

$$2\Phi^T (\mathbf{y} - \Phi\hat{\theta}) = 0$$

$$\Phi^T \mathbf{y} - \Phi^T \Phi \hat{\theta} = 0$$

$$\hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

# Parameter Estimation in Matrix form

- Least Squares Estimation

$$\hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \leftarrow \mathbf{y} = \Phi \theta + \epsilon, \quad \epsilon \sim N(0, \sigma^2 \mathbf{I})$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix} = \begin{bmatrix} \phi_1^T \\ \phi_2^T \\ \vdots \\ \phi_k^T \end{bmatrix} \theta + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_k \end{bmatrix}, \quad \Phi_k = \begin{bmatrix} \phi_{11} & \phi_{12} & \cdots & \phi_{1p} \\ \phi_{21} & \phi_{22} & \cdots & \phi_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{k1} & \phi_{k2} & \cdots & \phi_{kp} \end{bmatrix}$$

$$y_i = \phi_i^T \theta + \epsilon_i$$

$$y_i = \theta_0 + \theta_1 \phi_{i1} + \theta_2 \phi_{i2} + \cdots + \theta_p \phi_{ip} + \epsilon_i,$$

$$i = 1, \dots, k, \dots, n, \dots$$

- Observation Matrix

$$\Phi_k = [\phi_1 \ \phi_2 \ \cdots \ \phi_k]^T \xrightarrow{(p+1) \times k} \Phi_k^T \Phi_k = [\phi_1 \ \phi_2 \ \cdots \ \phi_k] \begin{bmatrix} \phi_1^T \\ \phi_2^T \\ \vdots \\ \phi_k^T \end{bmatrix} = \sum_{i=1}^k \phi_i \phi_i^T$$

$$\mathbf{y}_k = [y_1 \ y_2 \ \cdots \ y_k]^T$$

- Recursive Least Squares

$$\hat{\theta}_k = (\Phi_k^T \Phi_k)^{-1} \Phi_k^T \mathbf{y}_k \rightarrow \hat{\theta}_{k+1} = (\Phi_k^T \Phi_k + \phi_{k+1} \phi_{k+1}^T)^{-1} \Phi_{k+1}^T \mathbf{y}_{k+1}$$

# Parameter Estimation in Matrix form

- Matrix Inversion Lemma

$$(A + BDC)^{-1} = A^{-1} - A^{-1}B(D^{-1} + CA^{-1}B)^{-1}CA^{-1}$$

Sherman-Morrison formula:  $(A + uv^T)^{-1} = A^{-1} - \frac{A^{-1}uv^TA^{-1}}{1+v^TA^{-1}u}$

- Recursive Least Squares

$$\hat{\theta}_{k+1} = (\Phi_k^T \Phi_k + \phi_{k+1} \phi_{k+1}^T)^{-1} \Phi_{k+1}^T \mathbf{y}_{k+1}$$

define  $P_k \cong (\Phi_k^T \Phi_k)^{-1}$ ,

$$\begin{aligned} \hat{\theta}_{k+1} &= (P_k^{-1} + \phi_{k+1} \phi_{k+1}^T)^{-1} \Phi_{k+1}^T \mathbf{y}_{k+1} \\ &= \left( P_k - \frac{P_k \phi_{k+1} \phi_{k+1}^T P_k}{1 + \phi_{k+1}^T P_k \phi_{k+1}} \right) \Phi_{k+1}^T \mathbf{y}_{k+1}, \quad (\text{don't need inverse}) \end{aligned}$$

define  $G_k \cong \frac{P_k \phi_{k+1}}{1 + \phi_{k+1}^T P_k \phi_{k+1}} \Rightarrow P_{k+1} = P_k - \frac{P_k \phi_{k+1} \phi_{k+1}^T P_k}{1 + \phi_{k+1}^T P_k \phi_{k+1}} = P_k - G_k \phi_{k+1}^T P_k$

# Parameter Estimation in Matrix form

- Recursive Least Squares (cont.)

$$\begin{aligned}
 \hat{\theta}_{k+1} &= (P_k - G_k \phi_{k+1}^T P_k) [\Phi_k^T \quad \phi_{k+1}] \begin{bmatrix} \mathbf{y}_k \\ y_{k+1} \end{bmatrix} \\
 &= (P_k - G_k \phi_{k+1}^T P_k) (\Phi_k^T \mathbf{y}_k + \phi_{k+1} y_{k+1}) \\
 &= (I - G_k \phi_{k+1}^T) (P_k \Phi_k^T \mathbf{y}_k + P_k \phi_{k+1} y_{k+1}) \\
 &= (I - G_k \phi_{k+1}^T) (\hat{\theta}_k + P_k \phi_{k+1} y_{k+1}) \\
 &= \hat{\theta}_k - G_k \phi_{k+1}^T \hat{\theta}_k + P_k \phi_{k+1} y_{k+1} - G_k \phi_{k+1}^T P_k \phi_{k+1} y_{k+1} \\
 &= \hat{\theta}_k - G_k \phi_{k+1}^T \hat{\theta}_k + G_k y_{k+1} + G_k \phi_{k+1}^T P_k \phi_{k+1} y_{k+1} - G_k \phi_{k+1}^T P_k \phi_{k+1} y_{k+1}
 \end{aligned}$$

$$\hat{\theta}_{k+1} = \hat{\theta}_k + G_k (y_{k+1} - \phi_{k+1}^T \hat{\theta}_k), P_0 = \alpha \mathbf{I}, \alpha \gg 1.$$

$$\hat{\theta} = \hat{\theta}_n$$

$$G_k \cong \frac{P_k \phi_{k+1}}{1 + \phi_{k+1}^T P_k \phi_{k+1}}$$

$$P_{k+1} = P_k - G_k \phi_{k+1}^T P_k$$

$$\Phi_k = [\phi_1 \quad \phi_2 \quad \cdots \quad \phi_k]^T$$

$$\mathbf{y}_k = [y_1 \quad y_2 \quad \cdots \quad y_k]^T$$

$$P_k \cong (\Phi_k^T \Phi_k)^{-1}$$

$$G_k \cong \frac{P_k \phi_{k+1}}{1 + \phi_{k+1}^T P_k \phi_{k+1}}$$

$$P_{k+1} = P_k - G_k \phi_{k+1}^T P_k$$



# Parameter Estimation in Matrix form

- Weighted Recursive Least Squares

$$\hat{\theta}_{k+1} = (\lambda \Phi_k^T \Phi_k + \phi_{k+1} \phi_{k+1}^T)^{-1} \Phi_{k+1}^T \mathbf{y}_{k+1}, 0 < \lambda < 1$$

$$\hat{\theta}_{k+1} = \hat{\theta}_k + G_k (y_{k+1} - \phi_{k+1}^T \hat{\theta}_k), P_0 = \alpha \mathbf{I}, \alpha \gg 1$$

$$\hat{\theta} = \hat{\theta}_n$$

$$G_k \cong \frac{\lambda^{-1} P_k \phi_{k+1}}{1 + \lambda^{-1} \phi_{k+1}^T P_k \phi_{k+1}}$$

$$P_{k+1} = \lambda^{-1} P_k - \lambda^{-1} G_k \phi_{k+1}^T P_k$$

$$\Phi_k = [\phi_1 \ \phi_2 \ \cdots \ \phi_k]^T$$

$$\mathbf{y}_k = [y_1 \ y_2 \ \cdots \ y_k]^T$$

$$P_k \cong (\Phi_k^T \Phi_k)^{-1}$$

$$\lambda^{-1} P_k \cong (\lambda \Phi_k^T \Phi_k)^{-1}$$

# Quality of Fit in Matrix form

- Regression model in matrix form

$$\mathbf{y} = \Phi\theta + \epsilon$$

- Estimated parameter

$$\hat{\theta} = (\Phi^T\Phi)^{-1} \Phi^T \mathbf{y} = \theta + (\Phi^T\Phi)^{-1} \Phi^T \epsilon \text{ (unbiased estimate)}$$

- Confidence Interval

$$E(\hat{\theta}) = \theta,$$

$$\begin{aligned} E\left((\theta - \hat{\theta})^T(\theta - \hat{\theta})\right) &= E\epsilon^T \Phi (\Phi^T\Phi)^{-1} (\Phi^T\Phi)^{-1} \Phi^T \epsilon = E \text{Tr}(\epsilon^T \Phi (\Phi^T\Phi)^{-1} (\Phi^T\Phi)^{-1} \Phi^T \epsilon) \\ &= E \text{Tr}((\Phi^T\Phi)^{-1} (\cancel{\Phi^T\Phi})^{-1} \cancel{\Phi^T\Phi} \epsilon^T \epsilon) = \text{Tr}((\Phi^T\Phi)^{-1}) \sigma^2 \rightarrow \hat{\theta} = \theta \pm \alpha\sigma \end{aligned}$$

- Prediction

$$\hat{\mathbf{y}} = \Phi\hat{\theta} = \Phi\theta + \Phi(\Phi^T\Phi)^{-1} \Phi^T \epsilon = \Phi\theta + \mathbb{H}\epsilon,$$

where  $\mathbb{H}$  is symmetric and idempotent ( $\mathbb{H}^2 = \mathbb{H}$ ),  $\mathbb{H}\Phi = \Phi$ .

$$\mathbb{H}\hat{\mathbf{y}} = \mathbb{H}\Phi\theta + \mathbb{H}\epsilon = \Phi\theta + \mathbb{H}\epsilon = \hat{\mathbf{y}}$$

- Residual vector :  $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbb{H})\epsilon$

# Quality of Fit in Matrix form

- Residual vector

$$\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbb{H})\boldsymbol{\epsilon}$$

$$\begin{aligned} E(\mathbf{e}^T \mathbf{e}) &= E(\boldsymbol{\epsilon}^T (\mathbf{I} - \mathbb{H})(\mathbf{I} - \mathbb{H})\boldsymbol{\epsilon}) = E(\boldsymbol{\epsilon}^T (\mathbf{I} - \mathbb{H})\boldsymbol{\epsilon}) \\ &= \text{Tr}(\mathbf{I} - \mathbb{H})E(\boldsymbol{\epsilon}^T \boldsymbol{\epsilon}) = \text{Tr}(\mathbf{I} - \mathbb{H})n\sigma^2 = (n - p - 1)n\sigma^2 \end{aligned}$$

here

$$\begin{aligned} \text{Tr}(\mathbf{I} - \mathbb{H}) &= \text{Tr}(\mathbf{I}) - \text{Tr}(\mathbb{H}) = n - \text{Tr}(\Phi(\Phi^T \Phi)^{-1} \Phi^T) \\ &= n - \text{Tr}((\Phi^T \Phi)^{-1} \Phi^T \Phi) = n - (p + 1), p + 1: \# \text{ of parameters} \end{aligned}$$

hence

$$(p + 1) \times n \cdot n \times (p + 1)$$

$$E(\mathbf{e}^T \mathbf{e} / (n - p - 1)) = n\sigma^2 \rightarrow \frac{\mathbf{e}^T \mathbf{e}}{n - p - 1}: \text{unbiased estimate of } n\sigma^2$$

- Coefficient of Determination

$$R^2 = 1 - \frac{\mathbf{e}^T \mathbf{e}}{(\mathbf{y} - \bar{y}\mathbf{1})^T (\mathbf{y} - \bar{y}\mathbf{1})}, R_a^2 = 1 - \frac{\mathbf{e}^T \mathbf{e} / (n - p - 1)}{(\mathbf{y} - \bar{y}\mathbf{1})^T (\mathbf{y} - \bar{y}\mathbf{1}) / (n - 1)}$$

# PARTIAL LEAST SQUARES REGRESSION

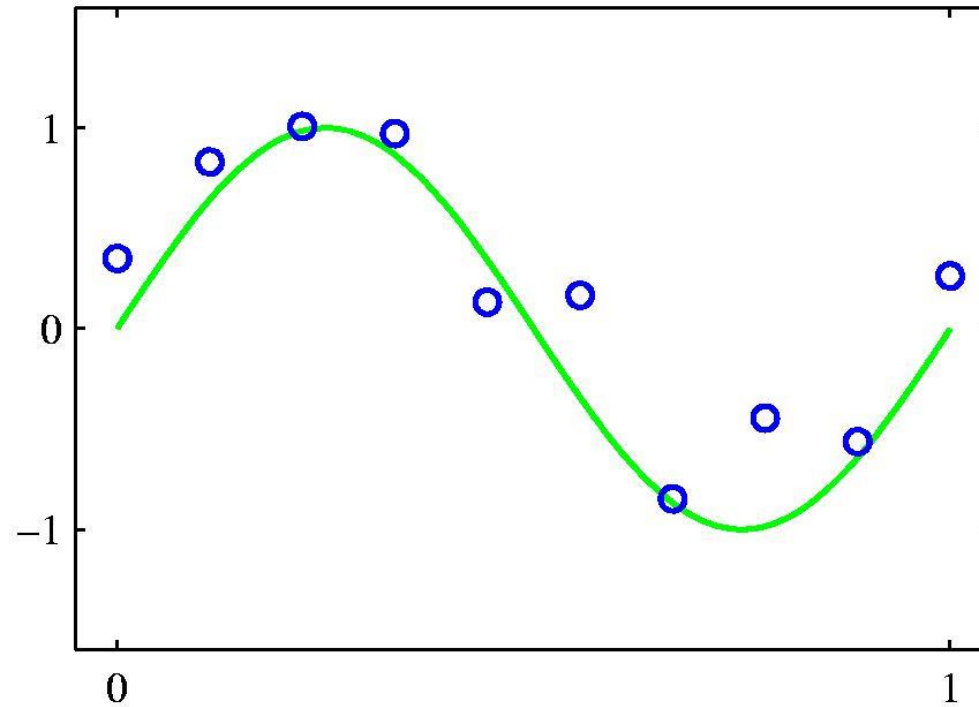
---

**JIN YOUNG CHOI**

**ECE, SEOUL NATIONAL UNIVERSITY**

# Overfitting and Underfitting

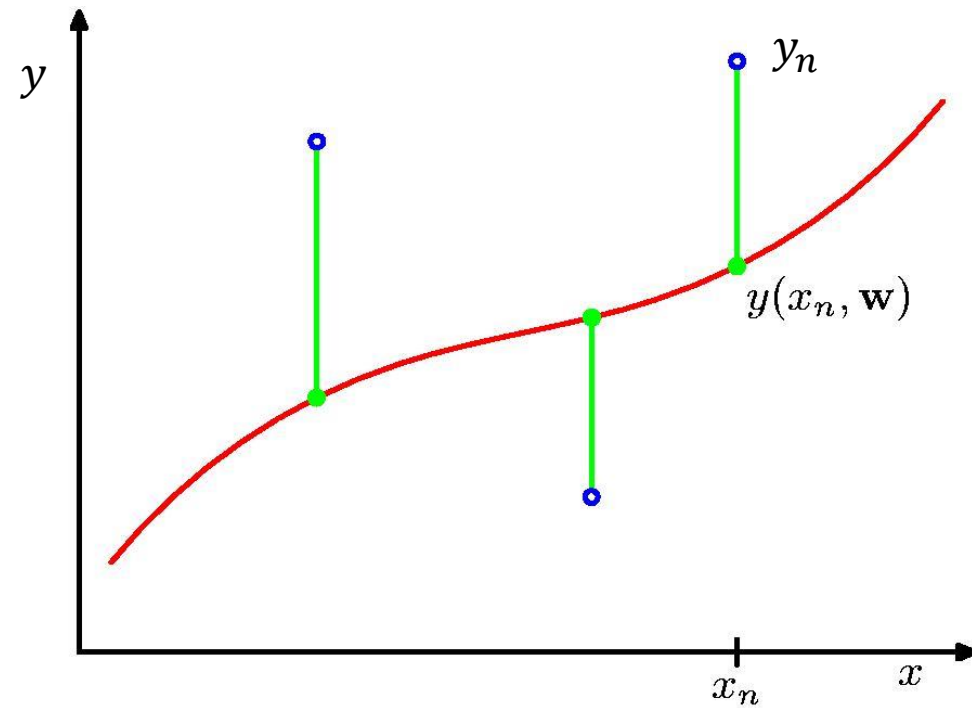
---



$$Y = \theta_0 + \theta_1 X + \theta_2 X^2 + \cdots + \theta_M X^M + \epsilon$$

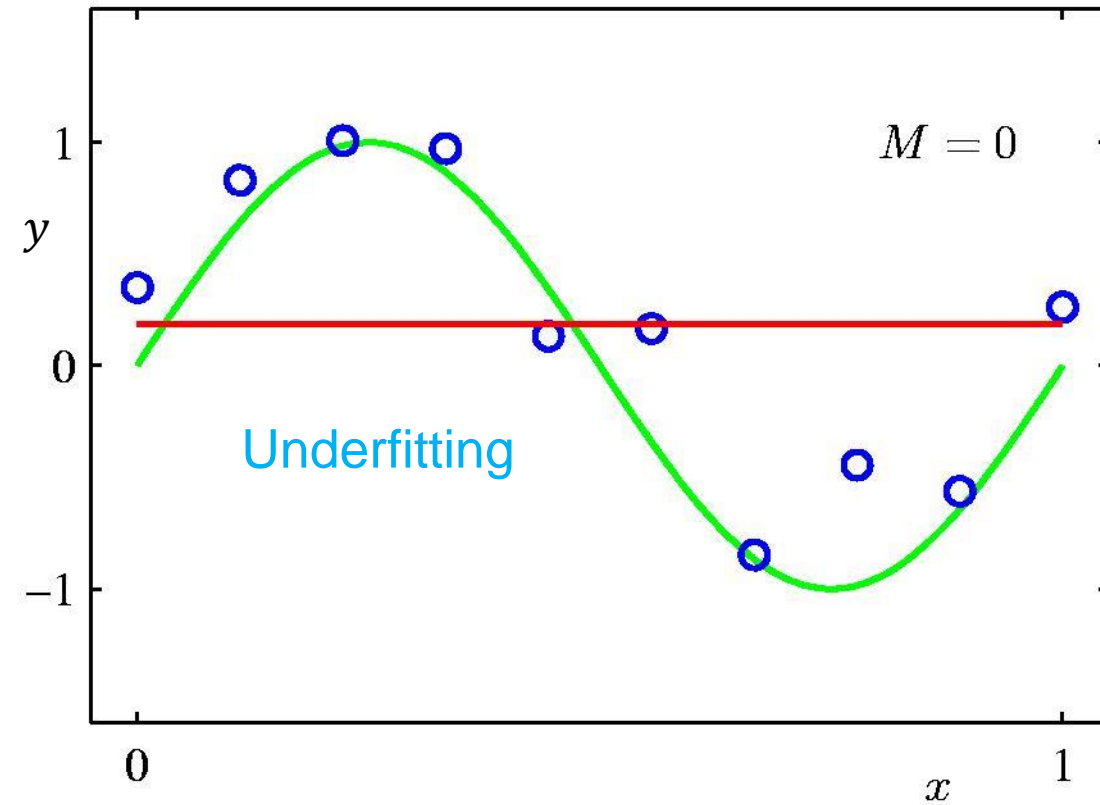
# Sum-of-Squares Error Function

---



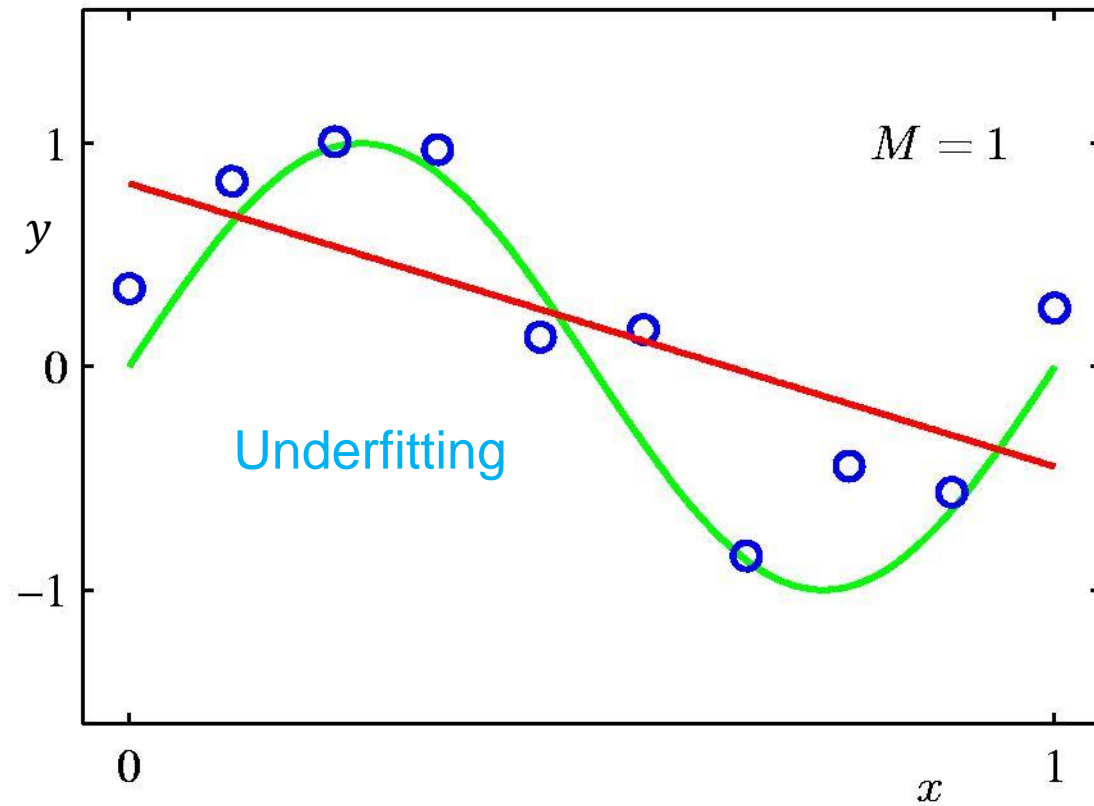
# 0<sup>th</sup> Order Polynomial

---



# 1<sup>st</sup> Order Polynomial

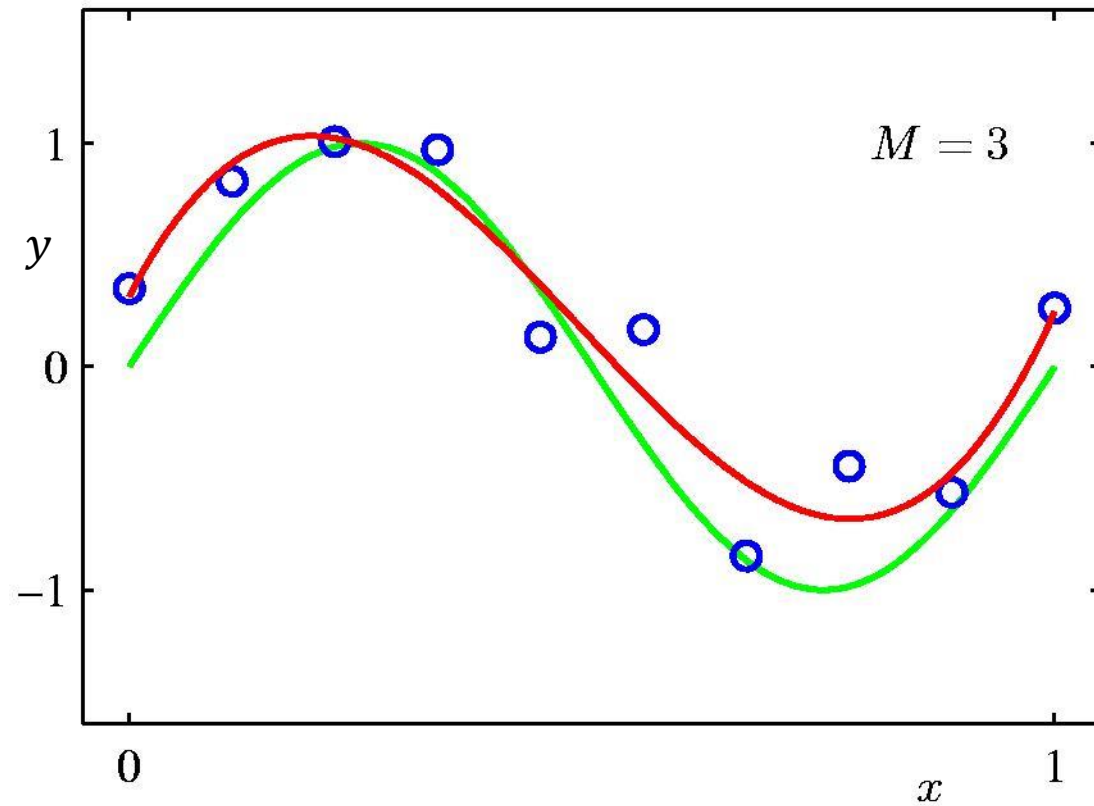
---





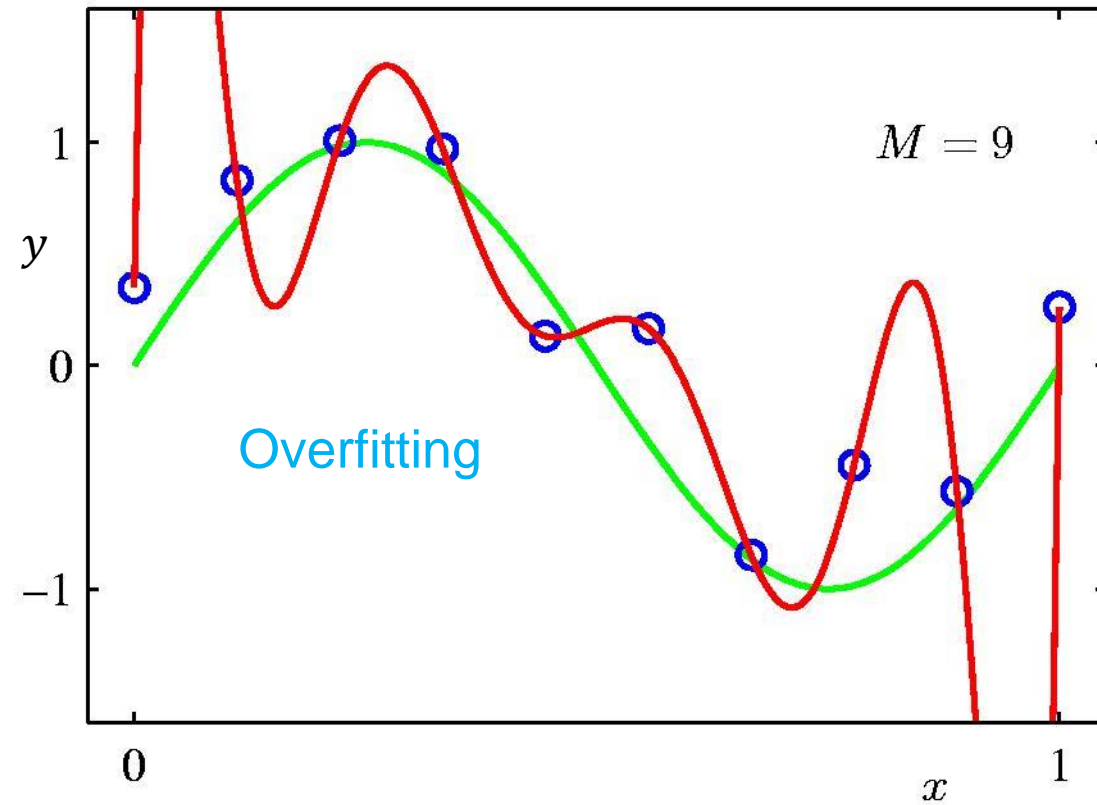
# 3<sup>rd</sup> Order Polynomial

---

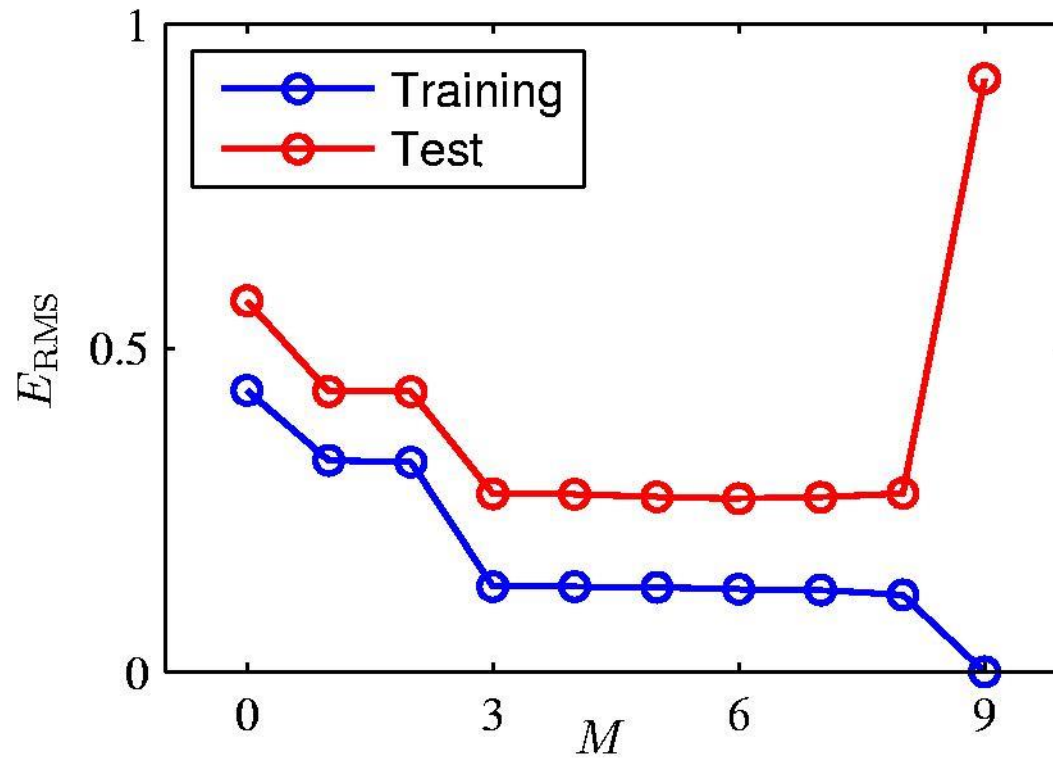


# 9<sup>th</sup> Order Polynomial

---



# Over-fitting

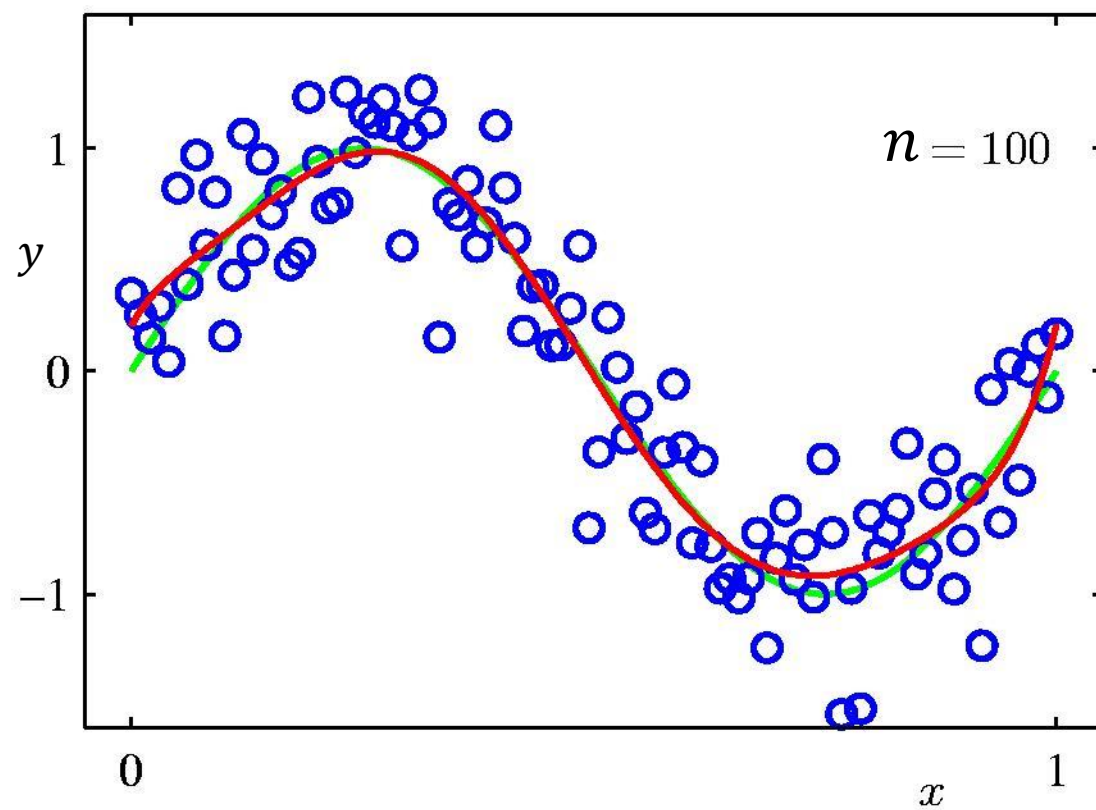


Root-Mean-Square (RMS) Error:  $E_{\text{RMS}} = \sqrt{E(\theta^*)/n}$

# Data Set Size:

---

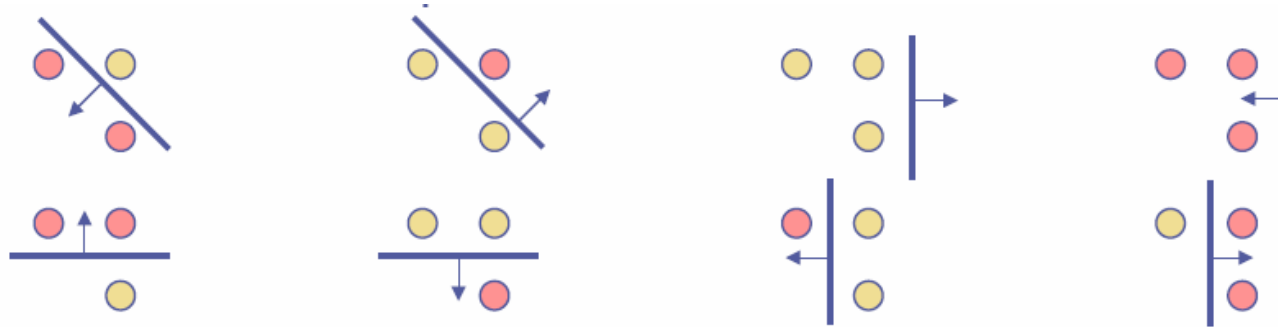
9<sup>th</sup> Order Polynomial



# Model Complexity

---

- VC(Vapnik-Chervonenkis )-dimension:  
Maximum number of points that can be labeled in all possible way
- VC dimension of linear classifiers in  $N$  dimensions  
is  $h=N+1$  (= #of weights,  $n_w$ ), cf.) MLP:  $O(n_w^2)$



- Measure of Complexity of a classifier
- **Minimizing VC dim. == Minimizing Complexity**

# Bias and Variance in Parameter Estimation

- Mean Squared Error(MSE) decomposition

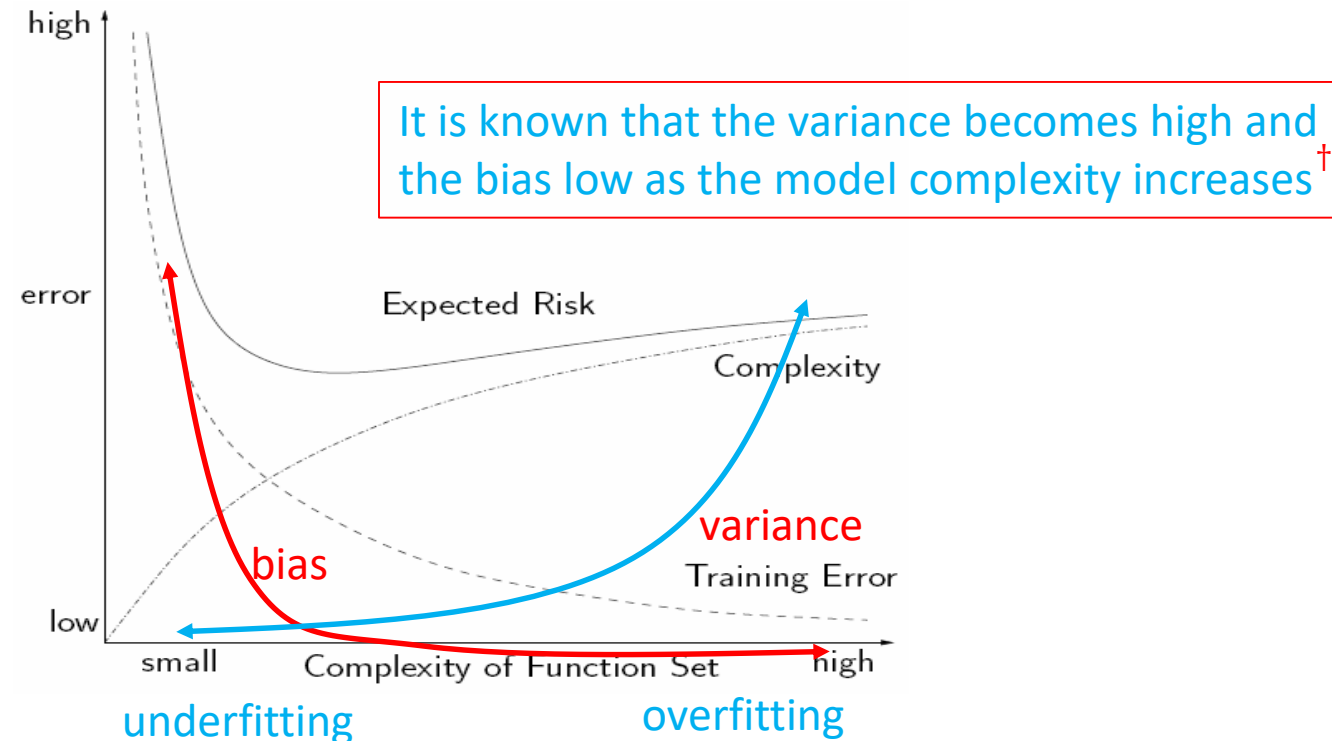
$$\begin{aligned}MSE(\hat{\theta}) &= E\left((\hat{\theta} - \theta)^2\right) \\&= E\left((\hat{\theta} - E(\hat{\theta}) + E(\hat{\theta}) - \theta)^2\right) \\&= E\left((\hat{\theta} - E(\hat{\theta}))^2 + 2(\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta) + (E(\hat{\theta}) - \theta)^2\right) \\&= E\left((\hat{\theta} - E(\hat{\theta}))^2\right) + 2 \overbrace{E\left((\hat{\theta} - E(\hat{\theta}))\right)}^{E(\hat{\theta}) - E(\hat{\theta}) = 0} (E(\hat{\theta}) - \theta) + (E(\hat{\theta}) - \theta)^2 \\&= E\left((\hat{\theta} - E(\hat{\theta}))^2\right) + (E(\hat{\theta}) - \theta)^2 \\&= Var(\hat{\theta}) + Bias(\hat{\theta}, \theta)^2\end{aligned}$$

overfitting      underfitting

# Bias and Variance in Parameter Estimation

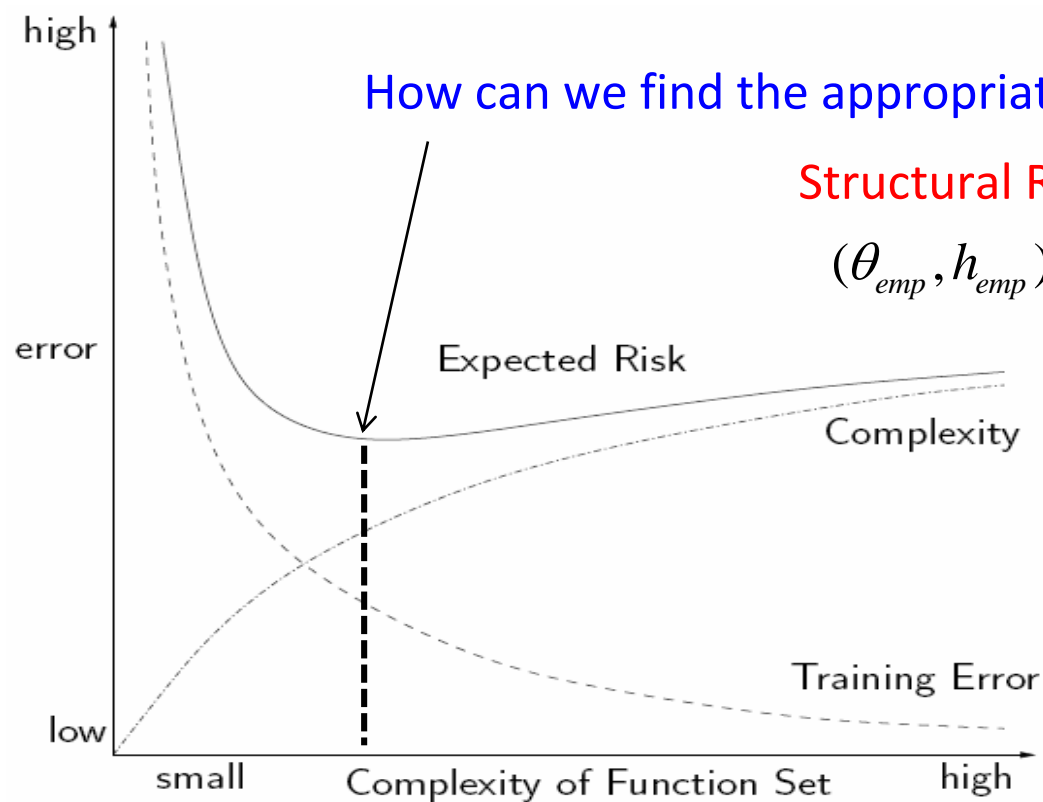
- Mean Squared Error(MSE) decomposition

$$\begin{aligned}MSE(\hat{\theta}) &= E((\hat{\theta} - \theta)^2) \\&= Var(\hat{\theta}) + Bias(\hat{\theta}, \theta)^2\end{aligned}$$



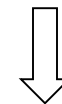
# Structural Risk Minimization

- For fixed training samples  $n$



Structural Risk Minimization

$$(\theta_{emp}, h_{emp}) = \arg \min_{\theta, h} R_{emp}(\theta, h)$$



Pruning (F-test, PCA)

Regularization

- ridge regression

- Lasso regression

SVM



# Partial Least Squares

- Matrix-vector form for General Regression (Revisit)

$$\text{Let } \theta = [\theta_0 \ \theta_1 \ \cdots \ \theta_M]^T, \ \phi_i = [1 \ \phi_{i1} \ \cdots \ \phi_{iM}]^T \\ \mathbf{y} = [y_1 \ y_2 \ \cdots \ y_n]^T, \ \boldsymbol{\epsilon} = [\epsilon_1 \ \epsilon_2 \ \cdots \ \epsilon_n]^T$$

$$\text{Then } y_i = \phi_i^T \theta + \epsilon_i, \ i = 1, \ \cdots, \ n.$$

$$\mathbf{y} = \Phi \theta + \boldsymbol{\epsilon}, \ \Phi = [\phi_1 \ \phi_2 \ \cdots \ \phi_n]^T$$

- Matrix-vector form for Multivariate Regression with no-intercept

$$y_i = \mathbf{x}_i^T \theta + \epsilon_i, \ i = 1, \ \cdots, \ n$$

$$\mathbf{y} = \mathbf{X} \theta + \boldsymbol{\epsilon}, \ \mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_n]^T$$

$$\mathbf{x}_i = [x_{i1} \ \cdots \ x_{ip}]^T, \ \theta = [\theta_1 \ \cdots \ \theta_p]^T \\ \mathbf{x}_i = \mathbf{x}_i^o - \mu, \ \mu = 1/n \sum_i \mathbf{x}_i^o$$

- Goal: reduce the input & parameter dimension:  $p > q$

$$\mathbf{x}_i = [x_{i1} \ \cdots \ x_{ip}]^T, \ \theta = [\theta_1 \ \cdots \ \theta_p]^T \longrightarrow \mathbf{z}_i = [z_{i1} \ \cdots \ z_{iq}]^T, \ \theta = [\theta_1 \ \cdots \ \theta_q]^T$$

# Principal Component Regression

$$\mathbf{a}_k = E^T(\mathbf{x}_k - \mathbf{m})$$

- Principal Component Analysis for  $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_n]$

$$\mathbf{S} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T = \mathbf{X}\mathbf{X}^T, \ \mathbf{S}\mathbf{u}_k = \lambda_k \mathbf{u}_k, \ \lambda_1 > \lambda_2 > \cdots > \lambda_p$$

$$\text{cov}(\mathbf{X}, \mathbf{X}) = \frac{1}{n} \mathbf{X} \mathbf{X}^T$$

- Reduced dim. vector ( $q < p$  dim.)

$$\mathbf{z}_i = \bar{\mathbf{U}}^T \mathbf{x}_i, \ \bar{\mathbf{U}} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_q]$$

$$q \times 1$$

$$p \times q$$

Orthonormal eigenvectors  $\{\mathbf{u}_i\}$

$$\mathbf{U}^T = \mathbf{U}^{-1}$$

$\mathbf{S}$  is symmetric

$$\mathbf{Z} = [\mathbf{z}_1 \ \mathbf{z}_2 \ \cdots \ \mathbf{z}_n] = \bar{\mathbf{U}}^T [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_n]$$

$$\mathbf{Z} = \bar{\mathbf{U}}^T \mathbf{X} \rightarrow \mathbf{Z}^T = \mathbf{X}^T \bar{\mathbf{U}}, \ \mathbf{y} = \mathbf{X}^T \boldsymbol{\theta} + \boldsymbol{\epsilon} \approx \mathbf{y} = \mathbf{Z}^T \bar{\mathbf{U}}^T \boldsymbol{\theta} + \boldsymbol{\epsilon} = \mathbf{y} = \mathbf{Z}^T \boldsymbol{\vartheta} + \boldsymbol{\epsilon}$$

- Applying LS algorithm to  $\mathbf{y} = \mathbf{Z}^T \boldsymbol{\vartheta} + \boldsymbol{\epsilon}$

$$\hat{\boldsymbol{\vartheta}} = \underset{\boldsymbol{\vartheta}}{\text{argmin}} \|\boldsymbol{\epsilon}\|^2 = \|\mathbf{y} - \mathbf{Z}^T \boldsymbol{\vartheta}\|^2 \rightarrow \hat{\boldsymbol{\vartheta}} = (\mathbf{Z}\mathbf{Z}^T)^{-1} \mathbf{Z}\mathbf{y} \rightarrow \hat{\mathbf{y}} = \mathbf{z}^T \hat{\boldsymbol{\vartheta}}, \ \mathbf{z} = \bar{\mathbf{U}}^T \mathbf{x}$$

# Partial Least Squares

- Nonlinear Iterative Partial Least Squares (NIPALS) algorithm

$$\mathbf{X}\mathbf{X}^T \mathbf{u} = \lambda \mathbf{u}$$

$$\text{Let } \mathbf{t} = \mathbf{X}^T \mathbf{u}$$

$$\mathbf{u} = \frac{1}{\lambda} \mathbf{X} \mathbf{t}$$

$$\text{Since } \|\mathbf{u}\| := 1 = \frac{1}{\lambda} \|\mathbf{X} \mathbf{t}\|$$

$$\lambda = \|\mathbf{X} \mathbf{t}\|$$

$$\bar{\mathbf{U}} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_q], \ \mathbf{z}_i = \bar{\mathbf{U}}^T \mathbf{x}_i$$

$$\mathbf{Z} = \bar{\mathbf{U}}^T \mathbf{X} \rightarrow \mathbf{Z}^T = \mathbf{X}^T \bar{\mathbf{U}}, \ \mathbf{Z} = [\mathbf{z}_1 \ \mathbf{z}_2 \ \cdots \ \mathbf{z}_n]$$

- Applying LS algorithm to  $\mathbf{y} = \mathbf{Z}^T \boldsymbol{\vartheta} + \boldsymbol{\epsilon}$

$$\hat{\boldsymbol{\vartheta}} = \underset{\boldsymbol{\vartheta}}{\operatorname{argmin}} \|\boldsymbol{\epsilon}\|^2 = \|\mathbf{y} - \mathbf{Z}^T \boldsymbol{\vartheta}\|^2 \rightarrow \hat{\boldsymbol{\vartheta}} = (\mathbf{Z}\mathbf{Z}^T)^{-1} \mathbf{Z} \mathbf{y} \rightarrow \hat{\mathbf{y}} = \mathbf{Z}^T \hat{\boldsymbol{\vartheta}}, \ \mathbf{z} = \bar{\mathbf{U}}^T \mathbf{x}$$

$\mathbf{t} := \mathbf{x}_j$  for some  $j$

Loop

$$\mathbf{u} = \mathbf{X} \mathbf{t} / \|\mathbf{X} \mathbf{t}\|$$

$$\mathbf{t} = \mathbf{X}^T \mathbf{u}$$

Until  $\mathbf{t}$  stop changing

$$\mathbf{X}^T := \mathbf{X}^T - \mathbf{t} \mathbf{u}^T = \mathbf{X}^T (\mathbf{I} - \mathbf{u} \mathbf{u}^T)$$

Repeat the Loop up to a small  $\|\mathbf{X} \mathbf{t}\|$

# Ridge Regression for Regularization

- $l_2$  regularization term is added

$$\hat{\theta} = \operatorname{argmin}_{\theta} \|\mathbf{y} - \Phi\theta\|^2 + \gamma \|\theta\|_2^2 (= S(\theta))$$

- solution:

$$\nabla_{\theta} ((\mathbf{y} - \Phi\theta)^T (\mathbf{y} - \Phi\theta) + \gamma \theta^T \theta) = 0 \text{ at } \hat{\theta}$$

$$2\Phi^T (\mathbf{y} - \Phi\hat{\theta}) + 2\gamma\hat{\theta} = 0$$

$$\hat{\theta} = (\Phi^T \Phi - \gamma \mathbf{I})^{-1} \Phi^T \mathbf{y}$$

$$\hat{\theta}_{k+1} = \hat{\theta}_k + G_k (y_{k+1} - \phi_{k+1}^T \hat{\theta}_k),$$

$$G_k \cong \frac{\lambda^{-1} P_k \phi_{k+1}}{1 + \lambda^{-1} \phi_{k+1}^T P_k \phi_{k+1}}$$

$$P_{k+1} = \lambda^{-1} P_k - \lambda^{-1} G_k \phi_{k+1}^T P_k, P_0 = -\gamma \mathbf{I}$$

$$P_0 = \alpha \mathbf{I}, \alpha \gg 1$$

# Lasso Regression for Regularization

---

- LASSO(Least Absolute Shrinkage Selector Operator)
- $l_1$  regularization term is added

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \|\mathbf{y} - \Phi\theta\|^2 + \gamma\|\theta\|_1$$

- solution:  $l_1$  norm is not differentiable  $\rightarrow$  constrained convex form by adding new optimization variables,

$$\begin{aligned} \hat{\theta} &= \underset{\theta}{\operatorname{argmin}} \|\mathbf{y} - \Phi\theta\|^2 + \gamma\mathbf{1}^T \mathbf{s} \\ &\text{subject to } |\theta_i| \leq s_i, \quad i = 1, \dots, n \end{aligned}$$

$$\begin{aligned} \hat{\theta} &= \underset{\theta}{\operatorname{argmin}} \|\mathbf{y} - \Phi\theta\|^2 + \gamma\mathbf{1}^T \mathbf{s} \\ &\text{subject to } -s_i \leq \theta_i \leq s_i, \quad i = 1, \dots, n \end{aligned}$$

# Elastic Regression for Regularization

---

- Ridge + LASSO

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \|\mathbf{y} - \Phi\theta\|^2 + \gamma_1 \|\theta\|_2^2 + \gamma_2 \|\theta\|_1$$

- solution:  $l_1$  norm is not differentiable  $\rightarrow$  constrained convex form by adding new optimization variables,

$$\begin{aligned} \hat{\theta} &= \underset{\theta}{\operatorname{argmin}} \|\mathbf{y} - \Phi\theta\|^2 + \gamma_1 \|\theta\|_2^2 + \gamma_2 \mathbf{1}^T \mathbf{s} \\ &\text{subject to } |\theta_i| \leq s_i, \quad i = 1, \dots, n \end{aligned}$$

$$\begin{aligned} \hat{\theta} &= \underset{\theta}{\operatorname{argmin}} \|\mathbf{y} - \Phi\theta\|^2 + \gamma_1 \|\theta\|_2^2 + \gamma_2 \mathbf{1}^T \mathbf{s} \\ &\text{subject to } -s_i \leq \theta_i \leq s_i, \quad i = 1, \dots, n \end{aligned}$$

# Interim Summary

---

- linear regression
  - simple linear regression
  - multiple linear regression
- nonlinear regression
  - logistic regression
  - high-order regression
  - basis-function regression
- matrix form for regression
  - recursive least squares
- partial least squares
  - over-fitting and underfitting
  - bias/variance
  - principle component regression
  - partial least squares algorithm
  - ridge regression
  - lasso, elastic regression
- Gaussian process regression

# Outline

---

- linear regression
  - simple linear regression
  - multiple linear regression
- nonlinear regression
  - logistic regression
  - high-order regression
  - basis-function regression
- matrix form for regression
  - recursive least squares
- partial least squares
  - over-fitting and underfitting
  - bias/variance
  - principle component regression
  - partial least squares algorithm
  - ridge regression
  - lasso, elastic regression
- Gaussian process regression



# GAUSSIAN PROCESS REGRESSION

**JIN YOUNG CHOI**

**ECE, SEOUL NATIONAL UNIVERSITY**

---

<https://arxiv.org/pdf/2009.10862.pdf>

<https://github.com/jwangjie/Gaussian-Processes-Regression-Tutorial>

<http://mlg.eng.cam.ac.uk/tutorials/06/es.pdf>

<https://www.sciencedirect.com/science/article/abs/pii/S0022249617302158>

<http://www.gaussianprocess.org/gpml/chapters/RW.pdf>

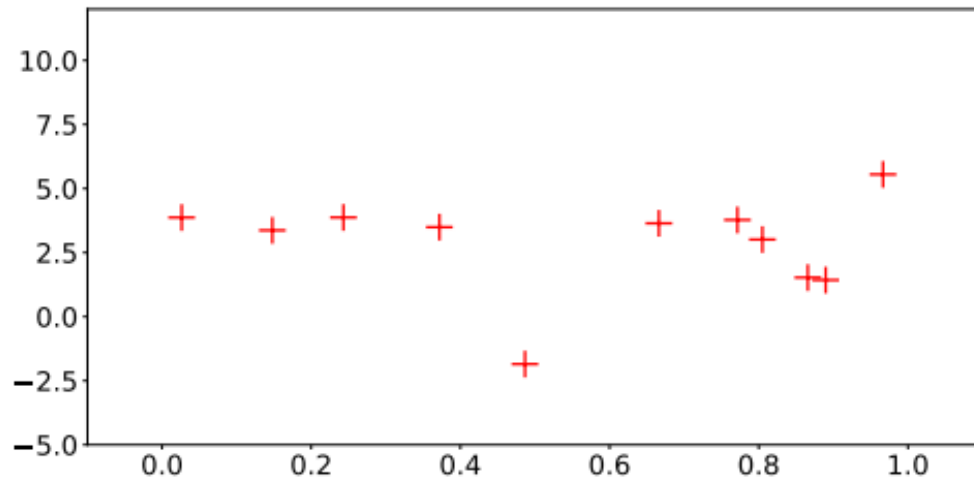
# Gaussian Process Regression

- General regression model (single variable)

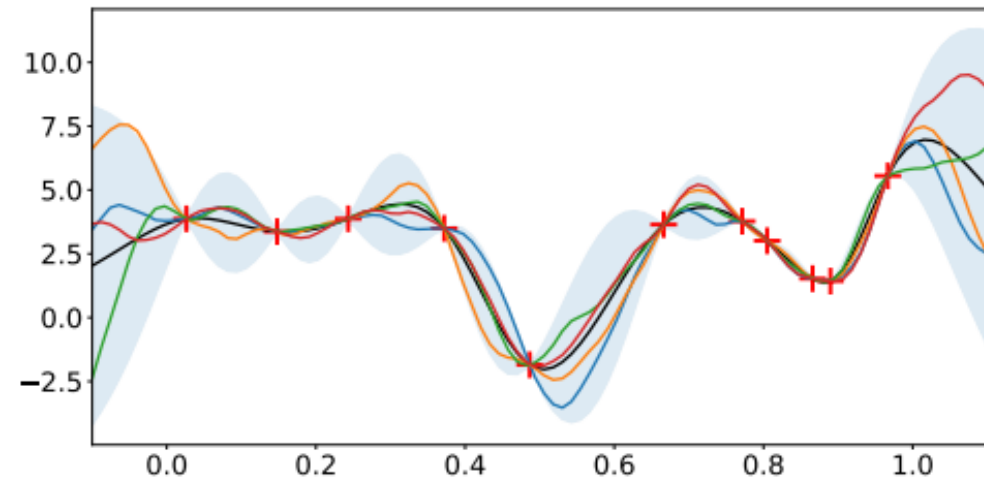
$$y = f(x) + \epsilon,$$

where  $\epsilon \sim N(0, \sigma^2)$  and so  $x, y$  are Gaussian random variables.

- Goal : to estimate  $f(x)$  with uncertainty from observation data  $D = \{(x_i, y_i) | i = 1, \dots, n\}$
- $x_i, y_i$  are treated as Gaussian random variables.



(a) Data point observations



(b) Five possible regression functions by GPR

# Gaussian Process Regression

---

- General regression model (single variable)

$$y = f(x) + \epsilon,$$

where  $\epsilon \sim N(0, \sigma^2)$  and so  $x, y$  are Gaussian random variables.

- Define

$$\mathbf{x}^T = [x_1 \quad \cdots \quad x_n], \quad \mathbf{y}^T = [y_1 \quad \cdots \quad y_n], \quad \mathbf{f} := \mathbf{f}(\mathbf{x}) = [f(x_1) \quad \cdots \quad f(x_n)].$$
$$p(\mathbf{x}, \mathbf{y}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \boldsymbol{\mu} \right)^T \Sigma^{-1} \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \boldsymbol{\mu} \right) \right] := \mathcal{N}(\boldsymbol{\mu}, \Sigma)$$

- Conditional probability (recall)

$$f_{X|Y}(x|y) = \frac{1}{(2\pi)^{\frac{k}{2}} \sqrt{\det \Sigma_{X|Y}}} \exp \left( -\frac{1}{2} (x - \mu_{X|Y})^t \Sigma_{X|Y}^{-1} (x - \mu_{X|Y}) \right),$$

where

$$\mu_{X|Y} = A(y - \mu_Y) + \mu_X \text{ and}$$

$$\Sigma_{X|Y} = \Sigma_X - AC_{YX}, \text{ where } A\Sigma_Y = \Sigma_{XY}.$$

# Gaussian Process Regression

- General regression model (single variable)

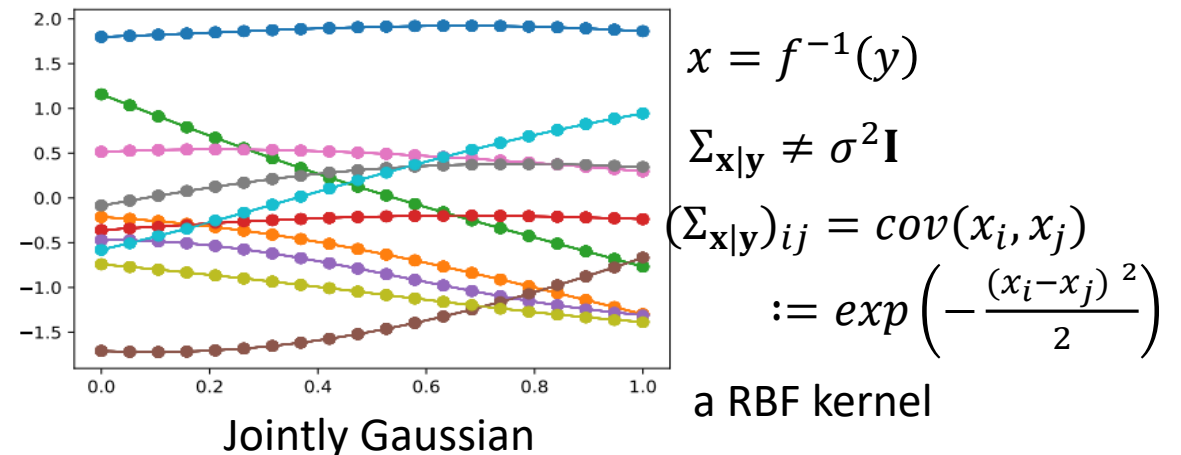
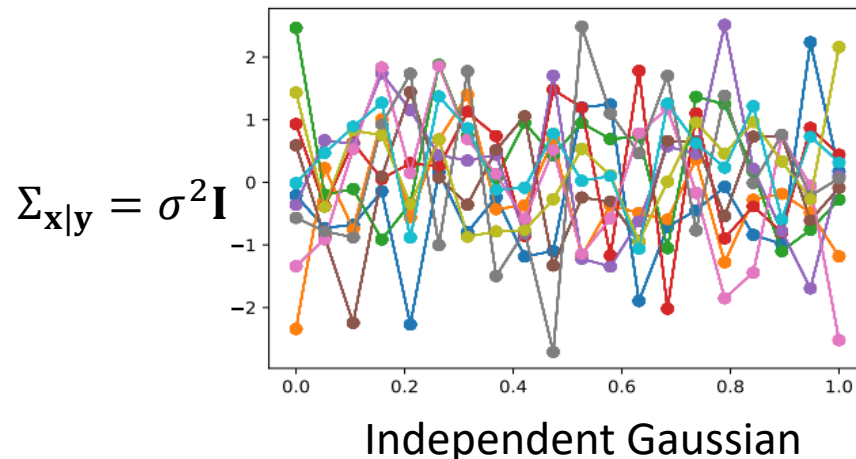
$$y = f(x) + \epsilon,$$

where  $\epsilon \sim N(0, \sigma^2)$  and so  $x, y$  are Gaussian random variables.

- Define

$$\mathbf{x} = [x_1 \quad \cdots \quad x_n], \quad \mathbf{y} = [y_1 \quad \cdots \quad y_n], \quad \mathbf{f} := \mathbf{f}(\mathbf{x}) = [f(x_1) \quad \cdots \quad f(x_n)].$$

$$p(\mathbf{x}|\mathbf{y}) = \frac{1}{(2\pi)^{d/2} |\Sigma_{\mathbf{x}|\mathbf{y}}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \mu_{\mathbf{x}|\mathbf{y}})^T \Sigma_{\mathbf{x}|\mathbf{y}}^{-1} (\mathbf{x} - \mu_{\mathbf{x}|\mathbf{y}}) \right] := \mathcal{N}(\mu_{\mathbf{x}|\mathbf{y}}, \Sigma_{\mathbf{x}|\mathbf{y}})$$



# Gaussian Process Regression

- Gaussian Processes ( $\mathcal{GP}$ ) for **multivariate** regression

$$y = f(\mathbf{x}) + \epsilon.$$

- define  $\mu_f(\mathbf{x}) := \mathbb{E}(f(\mathbf{x}))$ , then we assume  $f(\mathbf{x})$  is distributed as a Gaussian process

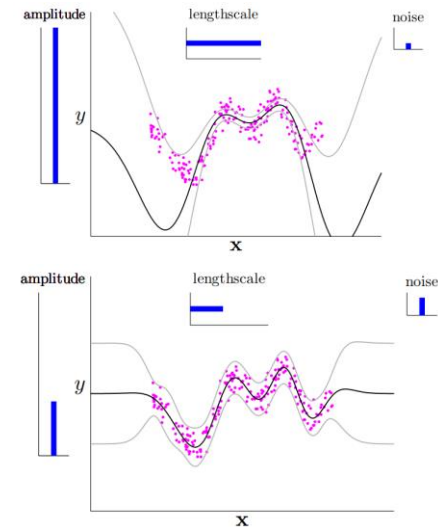
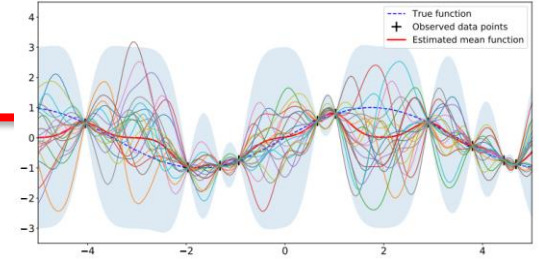
$$f(\mathbf{x}) \sim \mathcal{GP}(\mu_f(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

where  $k(\mathbf{x}, \mathbf{x}') = \mathbb{E} \left[ \left( f(\mathbf{x}) - \mu_f(\mathbf{x}) \right) \left( f(\mathbf{x}') - \mu_f(\mathbf{x}') \right) \right]$  called the kernel of  $\mathcal{GP}$ .

- The kernel is based on **assumptions** such as smoothness, that is, **similar  $\mathbf{x}, \mathbf{x}'$  yields similar  $f(\mathbf{x})$  and  $f(\mathbf{x}')$** . Thus a popular kernel is

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp \left( -\frac{1}{2\lambda} (\mathbf{x} - \mathbf{x}')^T (\mathbf{x} - \mathbf{x}') \right),$$

where hyperparameters  $\lambda$  and  $\sigma_f^2$  represents the length-scale and signal ( $f$ ) variance to control relation between  $\mathbf{x}$  and  $f(\mathbf{x})$ .



# Gaussian Process Regression

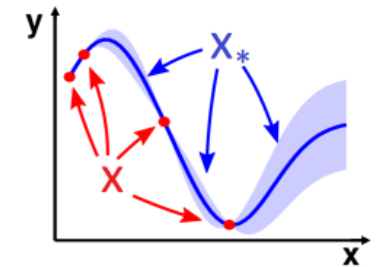
$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp \left( -\frac{1}{2\lambda} (\mathbf{x} - \mathbf{x}')^T (\mathbf{x} - \mathbf{x}') \right)$$

## Modeling of prior sampling function of $\mathcal{GP}$

- Denote  $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_n]$ ,  $\mathbf{y}^T = [y_1 \cdots y_n]$ ,  $\mathbf{f}^T := [f(\mathbf{x}_1) \cdots f(\mathbf{x}_n)]$ .

Let  $\mathbf{X}_*$  be a matrix containing a **new input points**  $\mathbf{x}_i^*$ ,  $i = 1, \dots, n$ . Then define the kernel matrix as

$$\mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) = \begin{bmatrix} k(\mathbf{x}_1^*, \mathbf{x}_1^*) & k(\mathbf{x}_1^*, \mathbf{x}_2^*) & \cdots & k(\mathbf{x}_1^*, \mathbf{x}_n^*) \\ k(\mathbf{x}_2^*, \mathbf{x}_1^*) & k(\mathbf{x}_2^*, \mathbf{x}_2^*) & \cdots & k(\mathbf{x}_2^*, \mathbf{x}_n^*) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n^*, \mathbf{x}_1^*) & k(\mathbf{x}_n^*, \mathbf{x}_2^*) & \cdots & k(\mathbf{x}_n^*, \mathbf{x}_n^*) \end{bmatrix}$$



- Choosing the prior mean function  $\mu_f(\mathbf{x}) = 0$ , we can sample values of  $f$  at inputs  $\mathbf{X}_*$  from  $\mathcal{GP}$  as

$$\mathbf{f}_* \sim \mathcal{N}(0, \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*))$$

which is the **prior distribution model** without observation data  $D = \{(x_i, y_i) | i = 1, \dots, n\}$ .

# Gaussian Process Regression

## Posterior predictions from a $\mathcal{GP}$

- Observations are  $D = \{(\mathbf{x}_i, y_i) | i = 1, \dots, n\} = \{\mathbf{X}, \mathbf{y}\}$ ,  $\mathbf{X} = [\mathbf{x}_1 \ \cdots \ \mathbf{x}_n]$ ,  $\mathbf{y}^T = [y_1 \ \cdots \ y_n]$ .
- The predictions for new inputs  $\mathbf{X}_*$  by drawing  $\mathbf{f}_*$  from the **posterior distribution**  $p(\mathbf{f} | D)$ .

A joint Gaussian distribution of  $\mathbf{y}$  and  $\mathbf{f}_*$ . Let  $\mathbf{X}_*$  follows

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_\epsilon^2 \mathbf{I} & \mathbf{K}(\mathbf{X}, \mathbf{X}_*) \\ \mathbf{K}(\mathbf{X}_*, \mathbf{X}) & \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right),$$

$$\begin{aligned} y &= f(x) + \epsilon \\ y_* &= f_*(x_*) + \epsilon \end{aligned}$$

where  $\sigma_\epsilon^2$  is the assumed noise level of the observations.

- The conditional distribution  $p(\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*)$  can be derived to a multivariate normal distribution with **mean**

$$\mu_{\mathbf{f}_*}(\mathbf{X}_*) = \mathbf{K}(\mathbf{X}_*, \mathbf{X})[\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_\epsilon^2 \mathbf{I}]^{-1} \mathbf{y}$$

**and variance**

$$\text{cov}_{\mathbf{f}_*}(\mathbf{X}_*) = \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) - \mathbf{K}(\mathbf{X}_*, \mathbf{X})[\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_\epsilon^2 \mathbf{I}]^{-1} \mathbf{K}(\mathbf{X}, \mathbf{X}_*)$$

# Gaussian Process Regression

## Posterior predictions from a $\mathcal{GP}$

- The mean function of the  $\mathcal{GP}$  can be given as

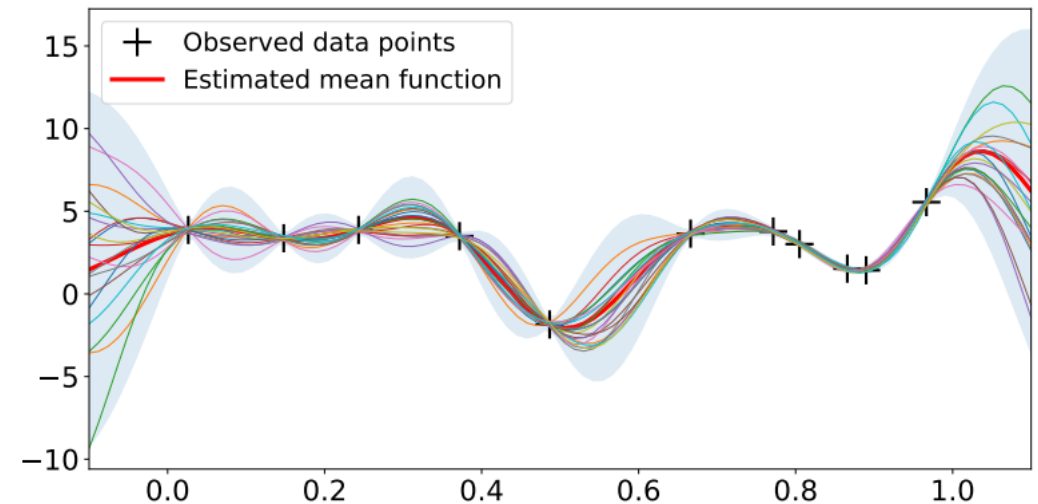
$$\mu_f(\mathbf{x}) = \mathbf{K}(\mathbf{x}, \mathbf{X})[\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_\epsilon^2 \mathbf{I}]^{-1} \mathbf{y}$$

and covariance function as

$$\text{cov}_f(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}') - \mathbf{K}(\mathbf{x}, \mathbf{X})[\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_\epsilon^2 \mathbf{I}]^{-1} \mathbf{K}(\mathbf{X}, \mathbf{x}')$$

$$\mathbf{K}(\mathbf{x}, \mathbf{X}) = [k(\mathbf{x}, \mathbf{x}_1) \quad k(\mathbf{x}, \mathbf{x}_2) \quad \cdots \quad k(\mathbf{x}, \mathbf{x}_n)]$$

$$\mathbf{K}(\mathbf{X}, \mathbf{x}') = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}') \\ k(\mathbf{x}_2, \mathbf{x}') \\ \vdots \\ k(\mathbf{x}_n, \mathbf{x}') \end{bmatrix}$$



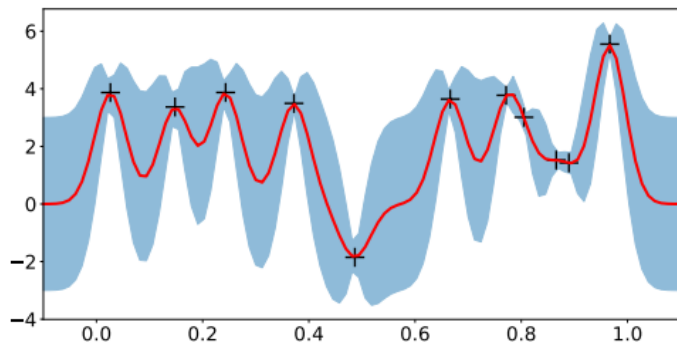


# Gaussian Process Regression

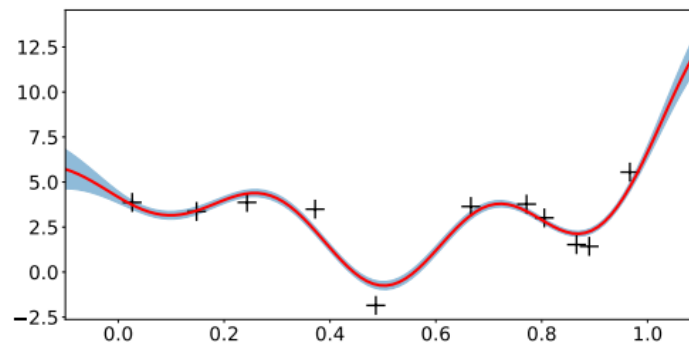
- The effect of the hyperparameters  $\lambda$  and  $\sigma_f^2$  of the kernel

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(-\frac{1}{2\lambda} (\mathbf{x} - \mathbf{x}')^T (\mathbf{x} - \mathbf{x}')\right) \approx \mathbb{E}\left[\left(f(\mathbf{x}) - \mu_f(\mathbf{x})\right)\left(f(\mathbf{x}') - \mu_f(\mathbf{x}')\right)\right],$$

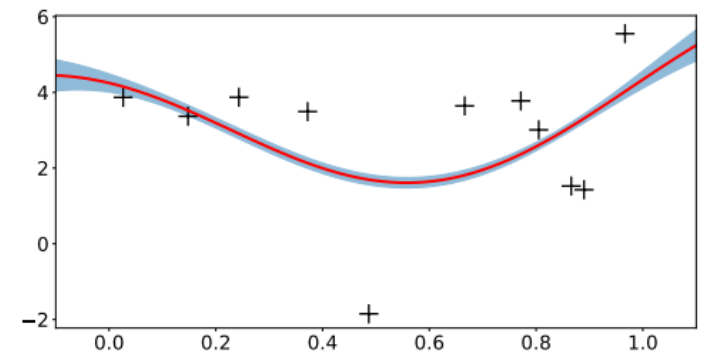
$\lambda$  : length-scale,  $\sigma_f^2$  : signal ( $f$ ) variance to control relation between  $\mathbf{x}$  and  $f(\mathbf{x})$ .



Small  $\lambda$



Medium  $\lambda$



Large  $\lambda$

# Gaussian Process Regression

$$p(\mathbf{y}|\mathbf{X}) = \frac{1}{(2\pi)^{d/2} |\Sigma_{\mathbf{y}|\mathbf{X}}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{y} - \mu_{\mathbf{y}|\mathbf{X}})^T \Sigma_{\mathbf{y}|\mathbf{X}}^{-1} (\mathbf{y} - \mu_{\mathbf{y}|\mathbf{X}}) \right]$$

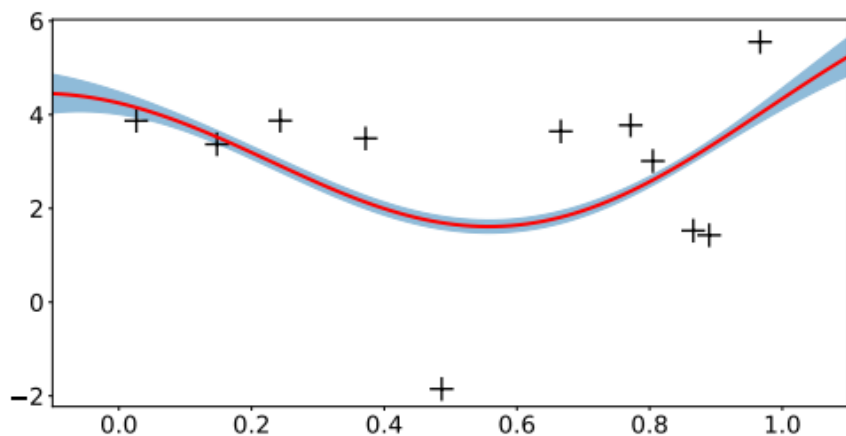
$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_\epsilon^2 \mathbf{I} & \mathbf{K}(\mathbf{X}, \mathbf{X}_*) \\ \mathbf{K}(\mathbf{X}_*, \mathbf{X}) & \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right)$$

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp \left( -\frac{1}{2\lambda} (\mathbf{x} - \mathbf{x}')^T (\mathbf{x} - \mathbf{x}') \right)$$

- The optimized hyperparameters  $\lambda$  and  $\sigma_f^2$

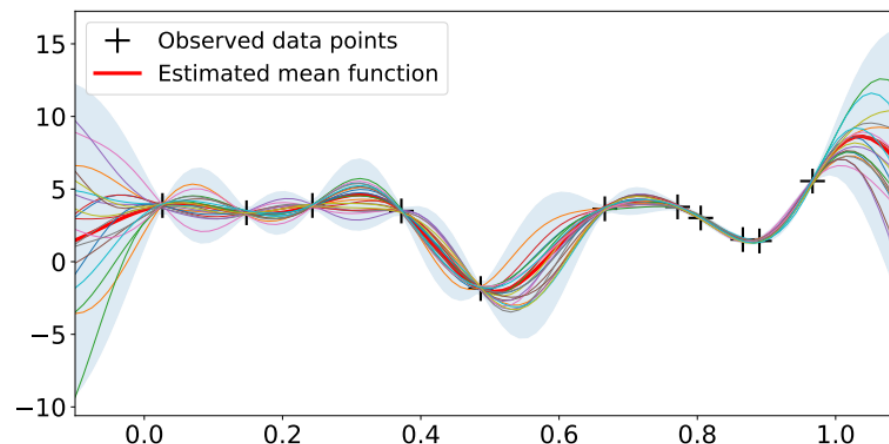
$$\lambda, \sigma_f^2 = \max_{\lambda, \sigma_f^2} \log p(\mathbf{y}|\mathbf{X})$$

$$\log p(\mathbf{y}|\mathbf{X}) = -\frac{1}{2} \mathbf{y}^T [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_\epsilon^2 \mathbf{I}]^{-1} \mathbf{y} - \frac{1}{2} \log \det[\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_\epsilon^2 \mathbf{I}] - \frac{n}{2} \log 2\pi$$



$$\sigma_f = 0.0067$$

$$\lambda = 0.0967$$



# Summary

---

- linear regression
  - simple linear regression
  - multiple linear regression
- nonlinear regression
  - logistic regression
  - high-order regression
  - basis-function regression
- matrix form for regression
  - recursive least squares
- partial least squares
  - over-fitting and underfitting
  - bias/variance
  - principle component regression
  - partial least squares algorithm
  - ridge regression
  - lasso, elastic regression
- Gaussian process regression