

# BLG 454E Learning From Data

FALL 2024-2025

Assoc. Prof. Yusuf Yaslan

Parametric Methods

# Parametric Estimation

- $\mathcal{X} = \{x^t\}_t$  where  $x^t \sim p(x)$
- Parametric estimation:  
Assume a form for  $p(x | \theta)$  and estimate  $\theta$ , its sufficient statistics, using  $X$   
e.g.,  $N(\mu, \sigma^2)$  where  $\theta = \{\mu, \sigma^2\}$

# Maximum Likelihood Estimation

- Likelihood of  $\theta$  given the sample  $\mathcal{X}$

$$l(\theta|\mathcal{X}) = p(\mathcal{X}|\theta) = \prod_t p(x^t|\theta)$$

- Log likelihood

$$L(\theta|\mathcal{X}) = \log l(\theta|\mathcal{X}) = \sum_t \log p(x^t|\theta)$$

- Maximum likelihood estimator (MLE)

$$\theta^* = \operatorname{argmax}_{\theta} L(\theta|\mathcal{X})$$

# Examples: Bernoulli/Multinomial

- **Bernoulli:** Two states, failure/success,  $x$  in  $\{0,1\}$

$$P(x) = p_o^x (1 - p_o)^{(1-x)}$$

$$\mathcal{L}(p_o | \mathcal{X}) = \log \prod_t p_o^{x^t} (1 - p_o)^{(1-x^t)}$$

$$\text{MLE: } p_o = \sum_t x^t / N$$

- **Multinomial:**  $K > 2$  states,  $x_i$  in  $\{0,1\}$

$$P(x_1, x_2, \dots, x_K) = \prod_i p_i^{x_i}$$

$$\mathcal{L}(p_1, p_2, \dots, p_K | \mathcal{X}) = \log \prod_t \prod_i p_i^{x_i^t}$$

$$\text{MLE: } p_i = \sum_t x_i^t / N$$

# Examples: Bernoulli (Derivation)

- **Bernoulli:** Two states, failure/success,  $x$  in  $\{0,1\}$

$$P(x) = p_o^x (1 - p_o)^{(1-x)}$$

$$L(p_o|X) = \log \prod_t p_o^{x^t} (1 - p_o)^{(1-x^t)}$$

$$\begin{aligned} \frac{dL(p_o | X)}{dp_o} &= \sum_{t=1}^N x^t \frac{d}{dp_o} \log(p_o) + \sum_{t=1}^N (1-x^t) \frac{d}{dp_o} \log(1-p_o) \\ &= \frac{1}{p_o} \sum_{t=1}^N x^t - \sum_{t=1}^N (1-x^t) \frac{1}{1-p_o} = 0 \end{aligned}$$

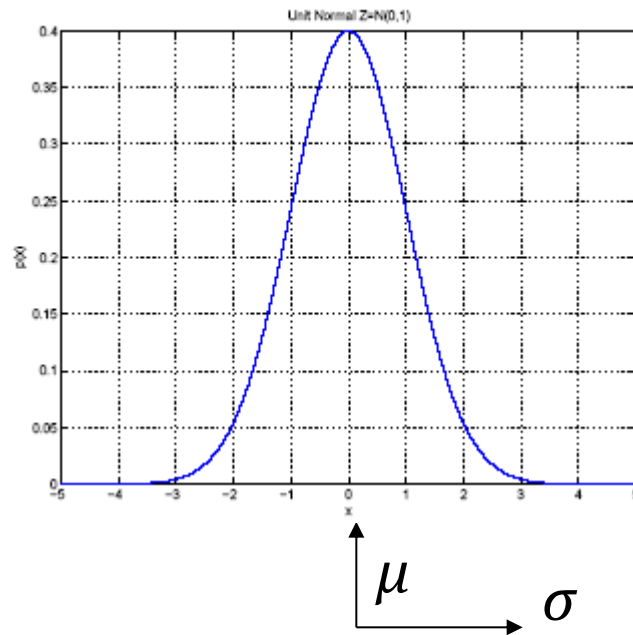
# Bernoulli (Derivation)

$$\square (1 - p_0) \sum_{t=1}^N x^t - p_0 \sum_{t=1}^N 1 - p_0 \sum_{t=1}^N x^t = 0$$

$$\square \sum_{t=1}^N x^t - p_0 N = 0 \Rightarrow p_0 = \frac{1}{N} \sum_{t=1}^N x^t$$

$$\text{MLE: } p_o = \sum_t x^t / N$$

# Gaussian (Normal) Distribution



- $p(x) = \mathcal{N}(\mu, \sigma^2)$

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

# Gaussian (Normal) Distribution

- Given that  $\mathcal{X} = \{x^t\}_t$  with  $x^t \sim \mathcal{N}(\mu, \sigma^2)$

$$L(\mu, \sigma | \mathcal{X}) = -\frac{N}{2} \log(2\pi) - N \log(\sigma) - \frac{\sum_{n=1}^N (x^t - \mu)^2}{2\sigma^2}$$

MLE for  $\mu$  and  $\sigma^2$ :

$$m \triangleq \frac{\sum_t x^t}{N}$$

$$s^2 \triangleq \frac{\sum_t (x^t - m)^2}{N}$$



# Probabilistic Interpretation of Linear Regression

$$r = f(x) + \varepsilon$$

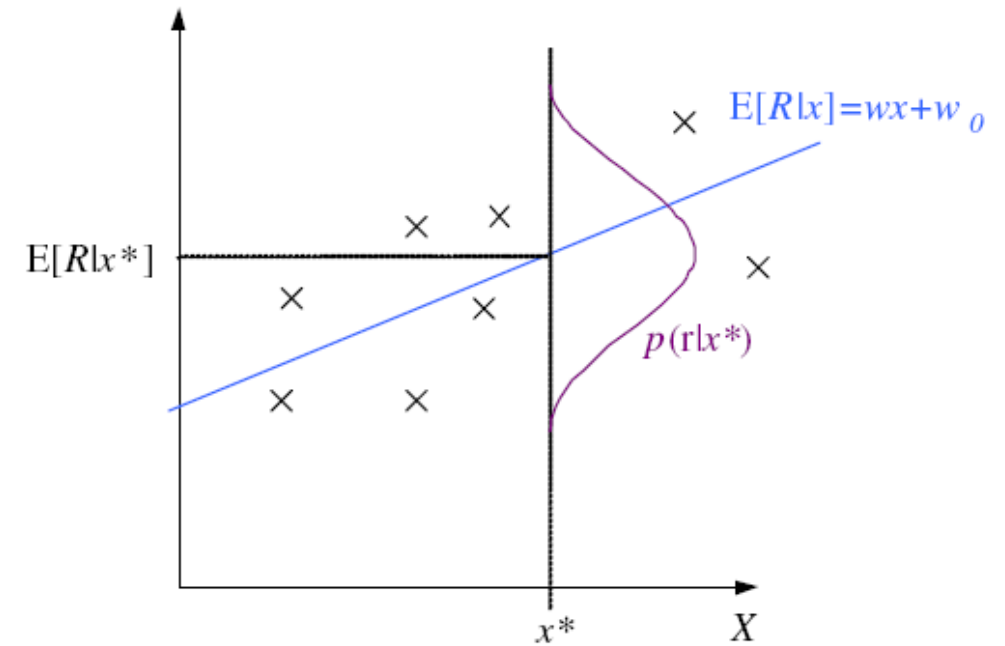
$$\text{estimator: } g(x | \theta)$$

$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$p(r | x) \sim \mathcal{N}(g(x | \theta), \sigma^2)$$

$$\mathcal{L}(\theta | \mathcal{X}) = -\log \prod_{t=1}^N p(x^t, r^t)$$

$$= -\log \prod_{t=1}^N p(r^t | x^t) - \log \prod_{t=1}^N p(x^t)$$



# Regression: From LogL to Error

$$\begin{aligned}\mathcal{L}(\theta|\mathcal{X}) &= \log \prod_{t=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{[r^t - g(x^t|\theta)]^2}{2\sigma^2} \right] \\ &= -N \log \sqrt{2\pi}\sigma - \frac{1}{2\sigma^2} \sum_{t=1}^N [r^t - g(x^t|\theta)]^2 \\ E(\theta|\mathcal{X}) &= \frac{1}{2} \sum_{t=1}^N [r^t - g(x^t|\theta)]^2\end{aligned}$$

# Linear Regression

$$E(\theta|\mathcal{X}) = \frac{1}{2} \sum_{t=1}^N [r^t - g(x^t|\theta)]^2$$

$$g(x^t | w_1, w_0) = w_1 x^t + w_0$$

Take derivative of E

$$\sum_t r^t = Nw_0 + w_1 \sum_t x^t$$

... wrto w0

$$\sum_t r^t x^t = w_0 \sum_t x^t + w_1 \sum_t (x^t)^2$$

... wrto w1

$$\mathbf{A} = \begin{bmatrix} N & \sum_t x^t \\ \sum_t x^t & \sum_t (x^t)^2 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} \sum_t r^t \\ \sum_t r^t x^t \end{bmatrix}$$

$$\mathbf{w} = \mathbf{A}^{-1} \mathbf{y}$$

# Polynomial Regression

$$g(x^t | w_k, \dots, w_2, w_1, w_0) = w_k x^t{}^k + \dots + w_2 x^t{}^2 + w_1 x^t + w_0$$

$$\mathbf{D} = \begin{bmatrix} 1 & x^1 & x^1{}^2 & \dots & x^1{}^k \\ 1 & x^2 & x^2{}^2 & \dots & x^2{}^k \\ \vdots & & & & \\ 1 & x^N & x^N{}^2 & \dots & x^N{}^k \end{bmatrix} \quad \mathbf{r} = \begin{bmatrix} r^1 \\ r^2 \\ \vdots \\ r^N \end{bmatrix}$$

$$\mathbf{w} = (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{D}^T \mathbf{r}$$

# Other Error Measures

- Square Error: 
$$E(\theta | \mathcal{X}) = \frac{1}{2} \sum_{t=1}^N (r^t - g(x^t | \theta))^2$$
- Relative Square Error: 
$$E(\theta | \mathcal{X}) = \frac{\sum_{t=1}^N (r^t - g(x^t | \theta))^2}{\sum_{t=1}^N (r^t - \bar{r})^2}$$
- Absolute Error:  $E(\vartheta | \mathcal{X}) = \sum_t |r^t - g(x^t | \vartheta)|$
- $\epsilon$ -sensitive Error: 
$$E(\vartheta | \mathcal{X}) = \sum_t 1(|r^t - g(x^t | \vartheta)| > \epsilon) (|r^t - g(x^t | \theta)| - \epsilon)$$

# Bias and Variance

Let  $X$  be a sample from a population specified up to a parameter  $\theta$

To evaluate the quality of this estimator we can measure how much it is different from  $\theta$

That is  $(d(X) - \theta)^2$

But since it is random variable (it depends on the sample) we need to average over all possible  $X$  and consider mean square error of the estimator

*Remember the properties of expectation*

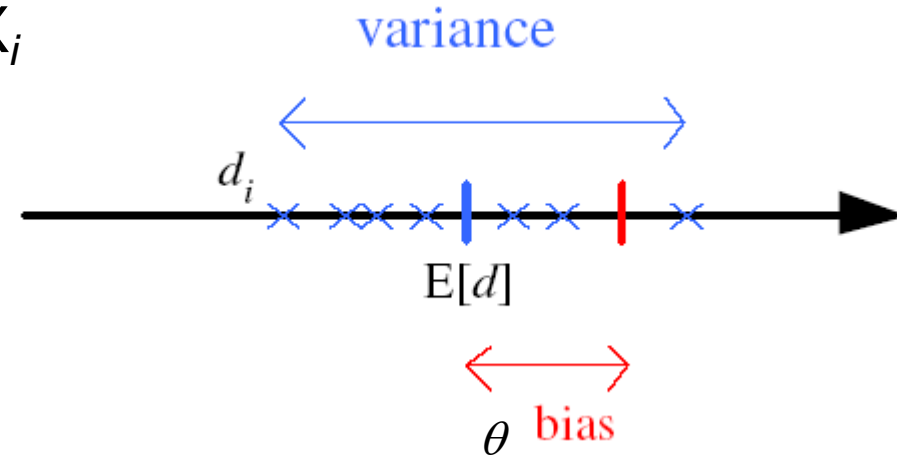
# Bias and Variance

Unknown parameter  $\theta$

Estimator  $d_i = d(X_i)$  on sample  $X_i$

Bias:  $b_{\theta}(d) = E[d] - \theta$

Variance:  $E[(d - E[d])^2]$



Mean square error:

$$r(d, \theta) = E[(d - \theta)^2] = E[(d - E[d] + E[d] - \theta)^2]$$

$$= (E[d] - \theta)^2 + E[(d - E[d])^2] + 2(d - E[d])(E[d] - \theta)$$

*Remember the properties of expectation*

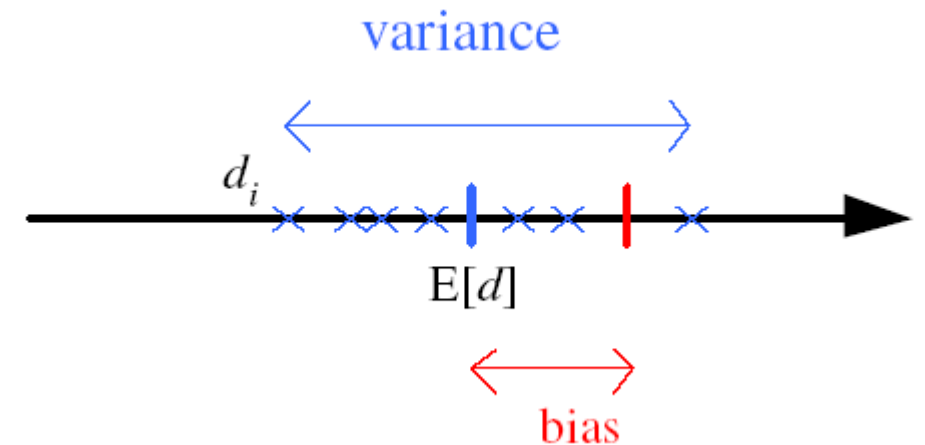
$$= E[(E[d]-\theta)^2] + E[(d-E[d])^2] + 2 E[(d-E[d])(E[d]-\theta)]$$

$$= E[(E[d]-\theta)^2] + E[(d-E[d])^2] + 2 (E[d]-E[d])(E[d]-\theta)$$

$$= (E[d] - \theta)^2 + E[(d-E[d])^2]$$

$$= (E[d] - \theta)^2 + E[(d-E[d])^2]$$

$$= \text{Bias}^2 + \text{Variance}$$





# Bias and Variance

$$E_x [E_x [r - g(x)]^2 | x] = E_x [r - E_x [r | x]]^2 + E_x [E_x [r | x] - g(x)]^2$$

*noise* *squared error*

$$E_x [E_x [r | x] - g(x)]^2 + E_x [E_x [r | x] - E_x [g(x)]]^2 = E_x [g(x) - E_x [g(x)]]^2$$

*bias* *variance*

# Estimating Bias and Variance

- $M$  samples  $X_i = \{x_i^t, r_i^t\}$ ,  $i=1, \dots, M$   
are used to fit  $g_i(x)$ ,  $i=1, \dots, M$  and  $t=1, \dots, N$

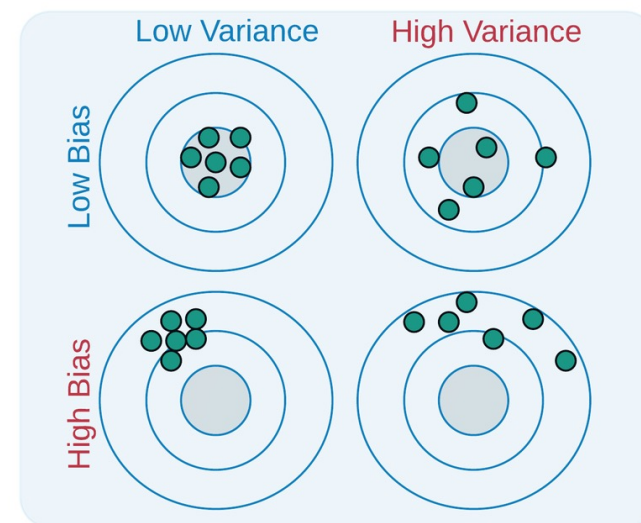
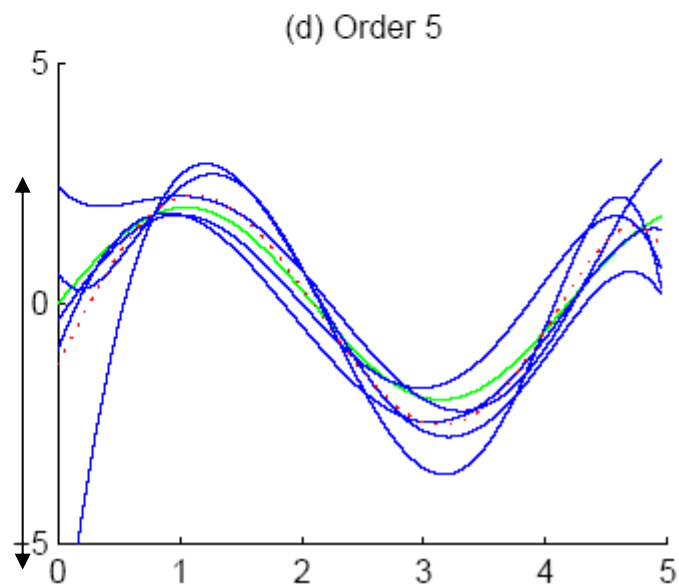
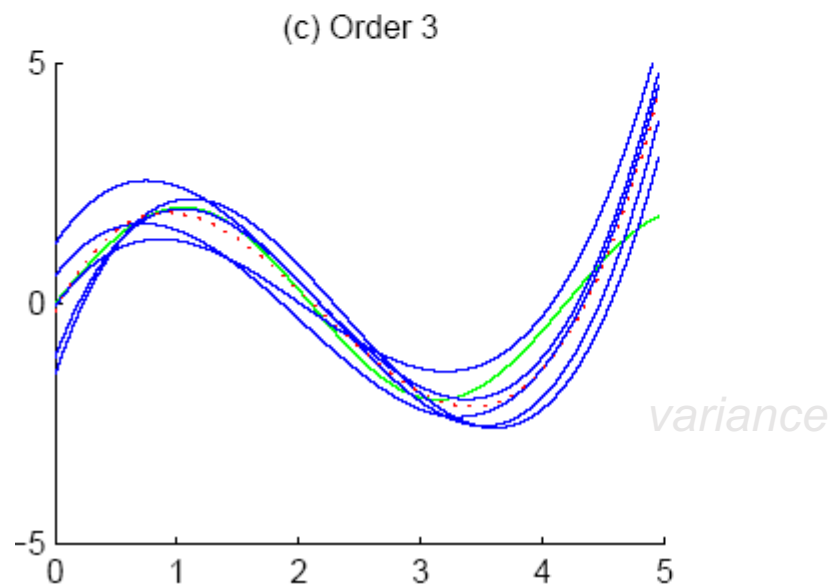
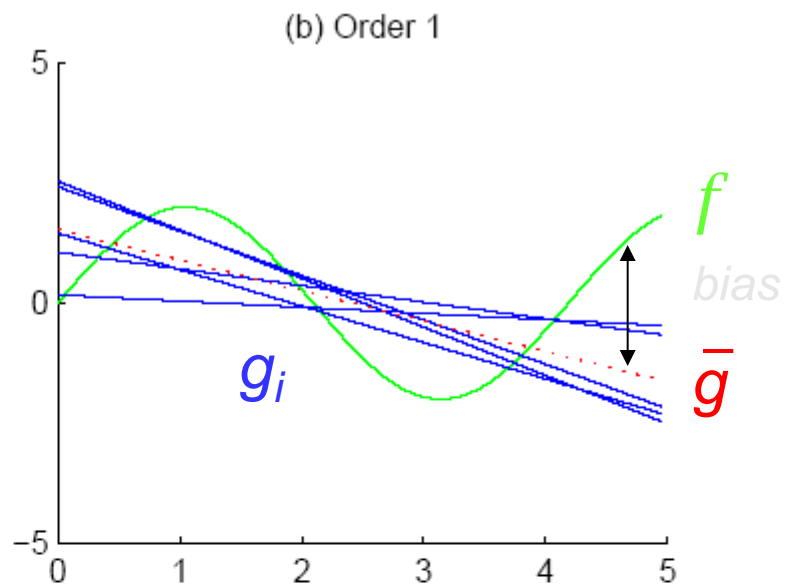
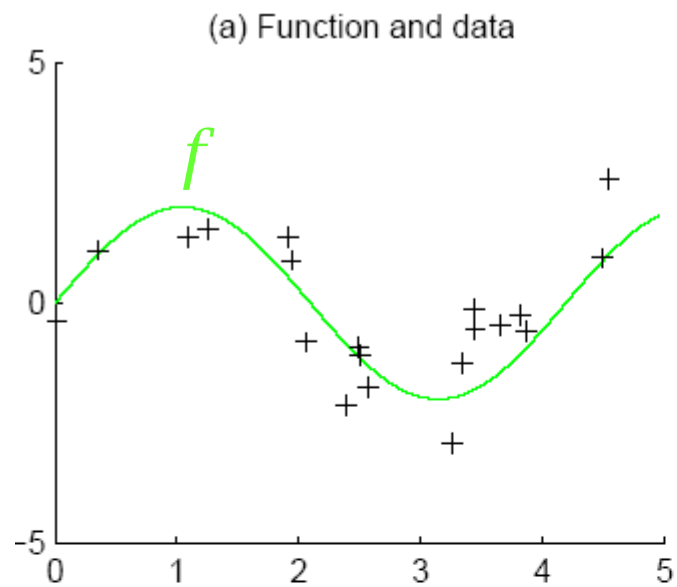
$$\text{Bias}^2[g] = \frac{1}{N} \sum_t [\bar{g}(x^t) - f(x^t)]^2$$

$$\text{Variance}[g] = \frac{1}{NM} \sum_t \sum_i [g_i(x^t) - \bar{g}(x^t)]^2$$

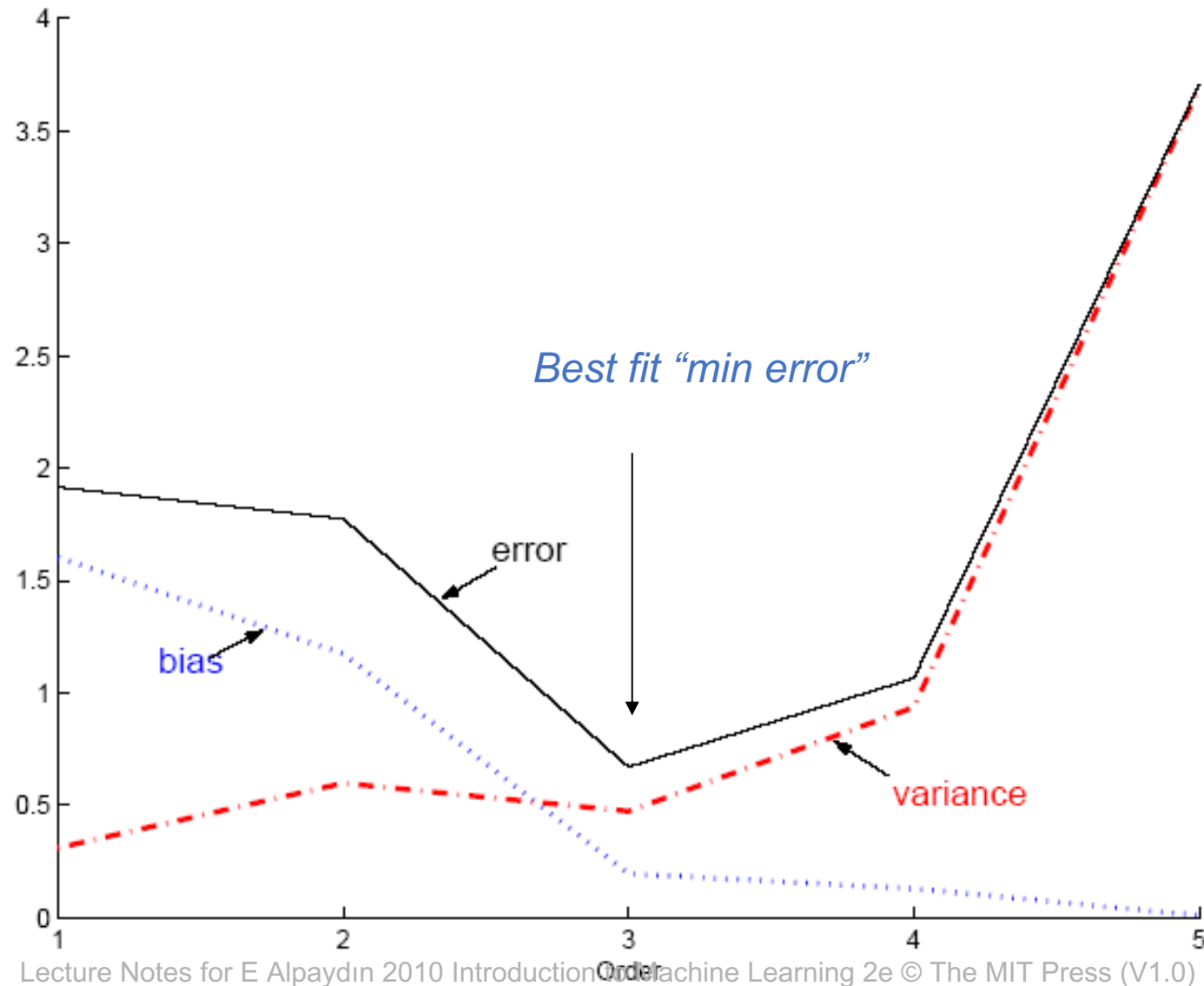
$$\bar{g}(x) = \frac{1}{M} \sum_i g_i(x)$$

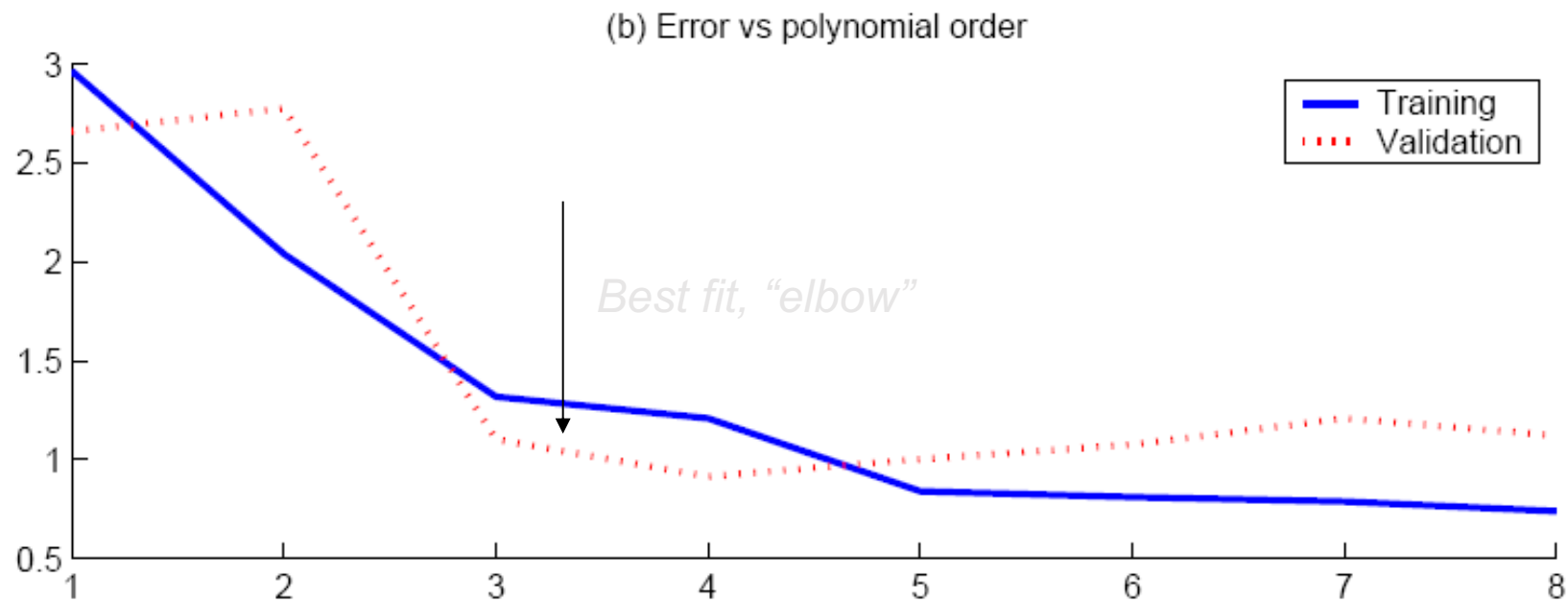
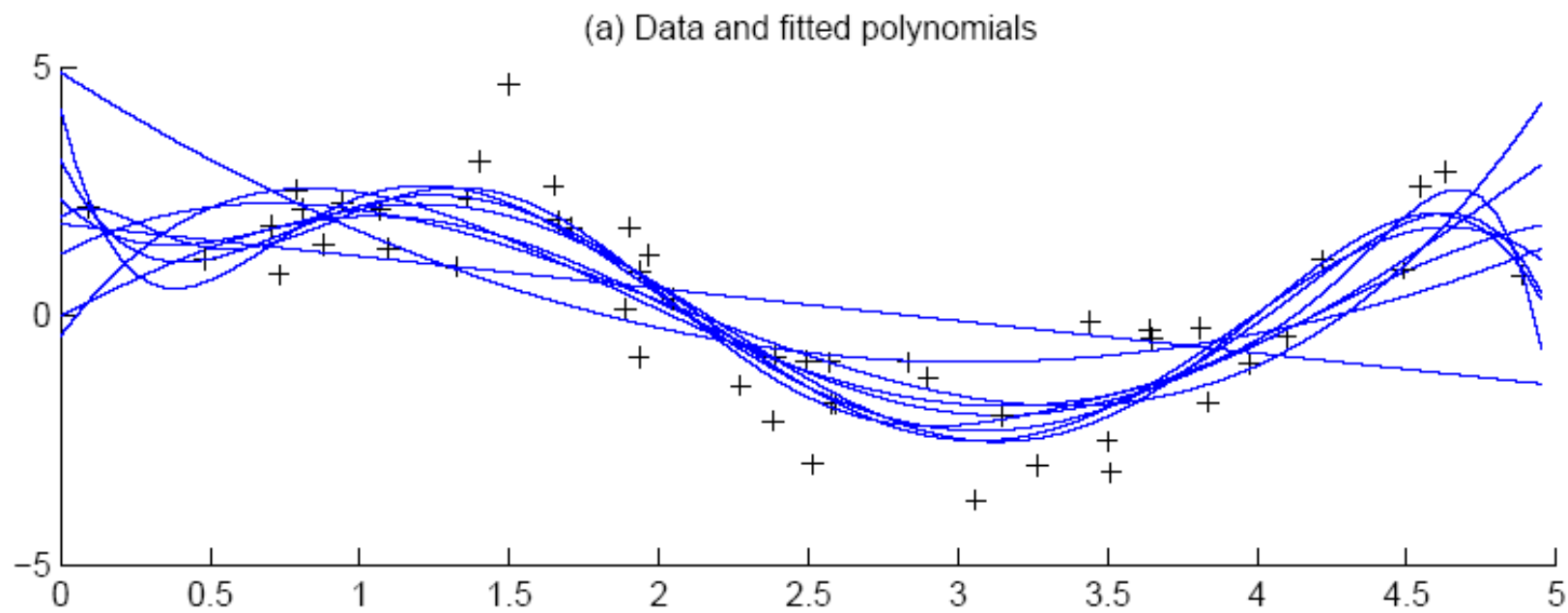
# Bias/Variance Dilemma

- Example:  $g_i(x)=2$  has no variance and high bias  
 $g_i(x)=\sum_t r_i^t/N$  has lower bias with variance
- As we increase complexity,  
    bias decreases (a better fit to data) and  
    variance increases (fit varies more with data)
- Bias/Variance dilemma: (Geman et al., 1992)

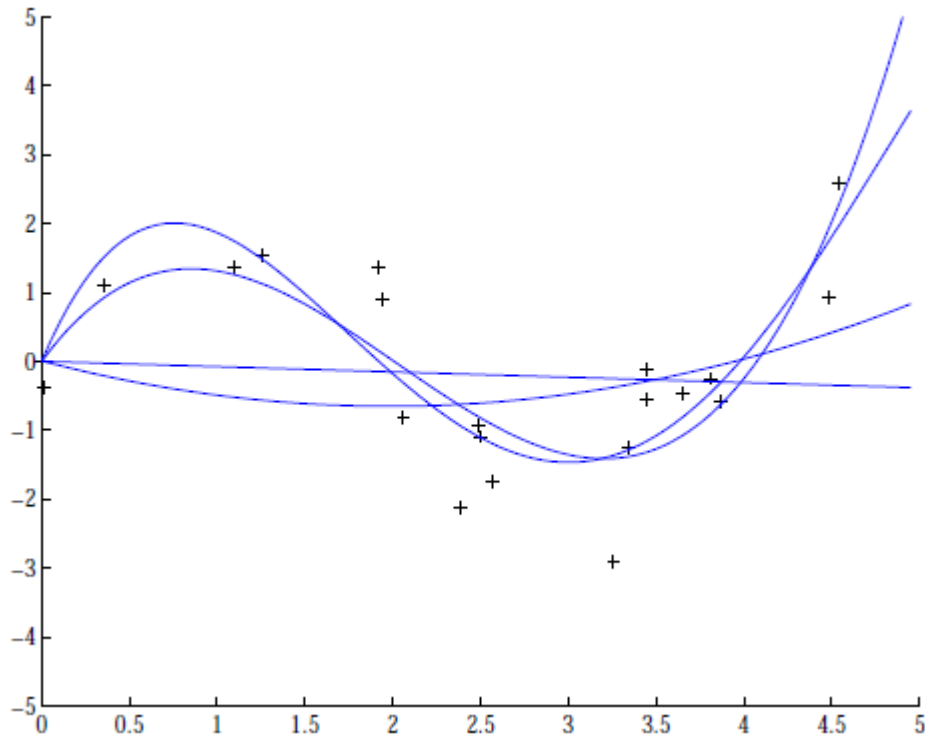


# Polynomial Regression





# Regression example



Coefficients increase in magnitude as order increases:

1:  $[-0.0769, 0.0016]$

2:  $[0.1682, -0.6657, 0.0080]$

3:  $[0.4238, -2.5778, 3.4675, -0.0002]$

4:  $[-0.1093, 1.4356, -5.5007, 6.0454, -0.0019]$

**Idea:** Penalize large coefficients

# Regularization

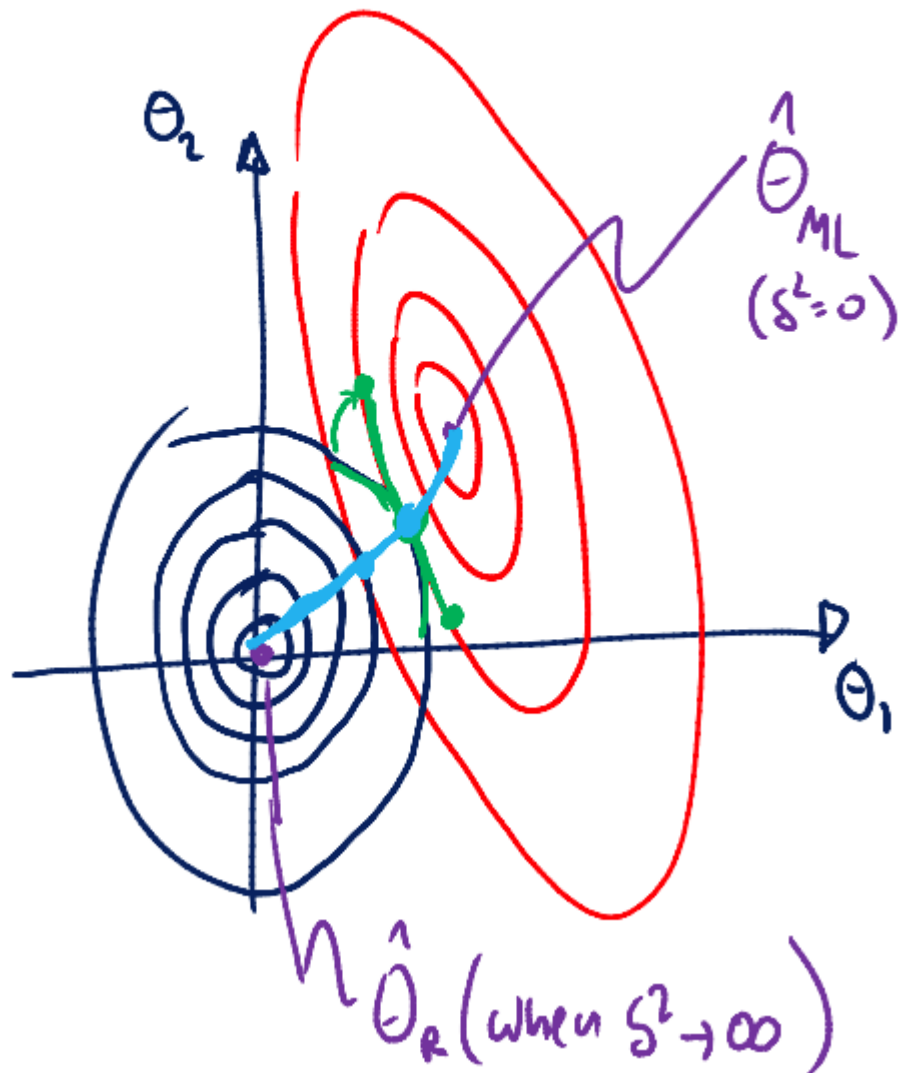
- New Cost Function  $E(\mathbf{w} | \mathcal{X}) = \frac{1}{2} \sum_{t=1}^N (y^t - g(x^t | \mathbf{w}))^2 + \lambda \sum_i w_i^2$
- Ridge Regression  $R(w) = \|\mathbf{w}\|^2 = \sum_i w_i^2$
- LASSO:  $R(w) = \|\mathbf{w}\|_1 = \sum_i |w_i|$

$$\mathcal{L}(W) = \frac{1}{2} \sum_{i=1}^N (y - Xw)^2 + \lambda \sum_i w_i^2 \Rightarrow \frac{1}{2} (y - Xw)^T (y - Xw) + \lambda w^T w$$



- $\nabla \mathcal{L} = -\frac{2}{2} X^T (y - Xw) + \lambda w$
- $-\frac{2}{2} X^T (y - Xw) + \lambda w = 0 \rightarrow X^T y = X^T X w + \lambda w \rightarrow$
- $X^T y = (X^T X + \lambda I) w$
- $\hat{w} = (X^T X + \lambda I)^{-1} X^T y$

$$J(\theta) = \underbrace{(y - X\theta)^T (y - X\theta)}_{\text{ellipses}} + \cancel{\delta^2 \theta^T \theta}$$



- Image is obtained from Nando Freitas' lecture notes

# Parametric Classification

$$g_i(x) = p(x|C_i)P(C_i)$$

or

$$g_i(x) = \log p(x|C_i) + \log P(C_i)$$

$$p(x|C_i) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{(x - \mu_i)^2}{2\sigma_i^2}\right]$$

$$g_i(x) = -\frac{1}{2}\log 2\pi - \log \sigma_i - \frac{(x - \mu_i)^2}{2\sigma_i^2} + \log P(C_i)$$

- Given the sample  $\mathcal{X} = \{x^t, r^t\}_{t=1}^N$   
 $x \in \mathfrak{R}$ 

$$r_i^t = \begin{cases} 1 & \text{if } x^t \in C_i \\ 0 & \text{if } x^t \in C_j, j \neq i \end{cases}$$

- ML estimates are

$$\hat{P}(C_i) = \frac{\sum_t r_i^t}{N} \quad m_i = \frac{\sum_t x^t r_i^t}{\sum_t r_i^t} \quad s_i^2 = \frac{\sum_t (x^t - m_i)^2 r_i^t}{\sum_t r_i^t}$$

- Discriminant becomes

$$g_i(x) = -\frac{1}{2} \log 2\pi - \log s_i - \frac{(x - m_i)^2}{2s_i^2} - \log \hat{P}(C_i)$$

