# BLG 454E Learning From Data

FALL 2022-2023

Assoc. Prof. Yusuf Yaslan

## Bias Variance

# Linear Regression

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

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

Take derivative of E

$$\sum_t r^t = N w_0 + w_1 \sum_t x^t$$

...wrto w0

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

...wrto w1

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

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

# Polynomial Regression

$$g\left(x^t \mid w_k, \ldots, w_2, w_1, w_0\right) = w_k\left(x^t\right)^k + \cdots + w_2\left(x^t\right)^2 + w_1 x^t + w_0$$

$$\mathbf{D} = \begin{bmatrix} 1 & x^1 & \left(x^1\right)^2 & \cdots & \left(x^1\right)^k \\ 1 & x^2 & \left(x^2\right)^2 & \cdots & \left(x^2\right)^k \\ \vdots & & & & \\ 1 & x^N & \left(x^N\right)^2 & \cdots & \left(x^N\right)^2 \end{bmatrix} \quad \mathbf{r} = \begin{bmatrix} r^1 \\ r^2 \\ \vdots \\ r^N \end{bmatrix}$$

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

# Other Error Measures

- Square Error:

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

- Relative Square Error:

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

- Absolute Error: $E(\vartheta \mid X) = \sum_t |r^t - g(x^t \mid \vartheta)|$

- ε-sensitive Error:

$$E(\vartheta \mid X) = \sum_t 1(|r^t - g(x^t \mid \vartheta)| > \varepsilon)\,(|r^t - g(x^t \mid \vartheta)| - \varepsilon)$$

# 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 meas square error of the estimator

*Remember the properties of expectation*

# Bias and Variance
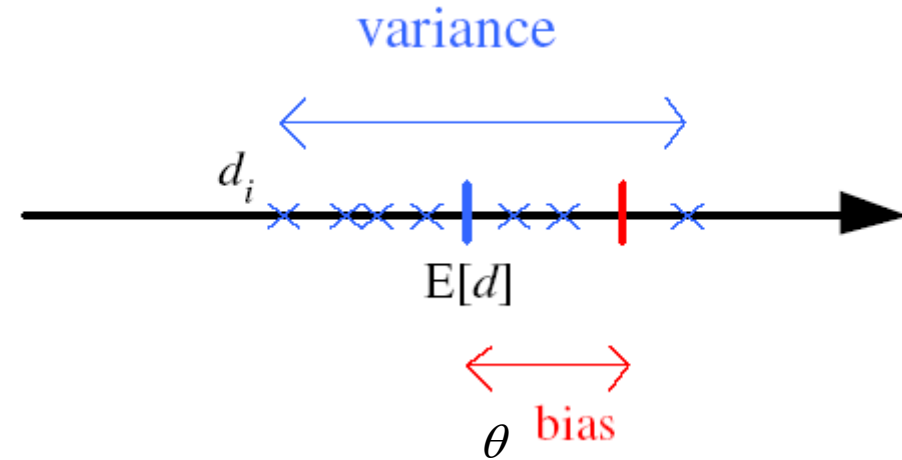
Unknown parameter $\theta$
Estimator $d_i = d(X_i)$ on sample $X_i$

Bias: $b_\theta(d) = E[d] - \theta$
Variance: $E[(d - E[d])^2]$

Mean square error:

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

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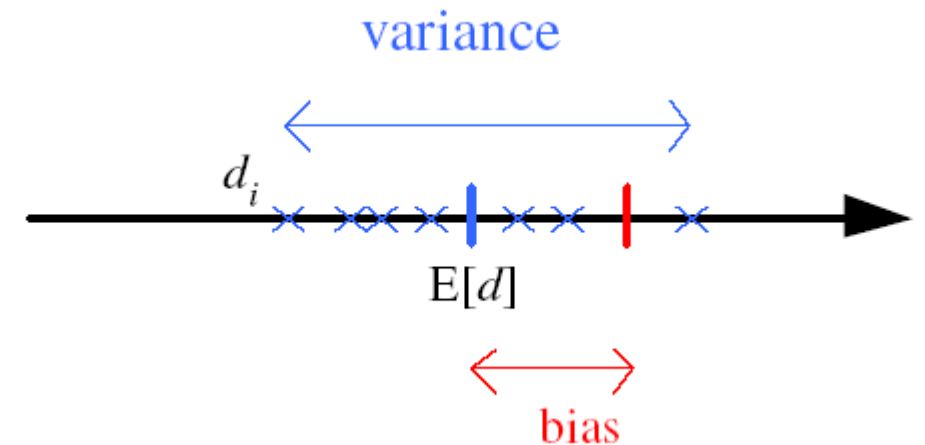


*Remember the properties of expectation*

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

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

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

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

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

# Bias and Variance

$$E\left[(r-g(x))^2 \mid x\right] = E\left[(r-E[r \mid x])^2 \mid x\right] + (E[r \mid x] - g(x))^2$$

*noise*         *squared error*

$$E_x\left[(E[r \mid x] - g(x))^2 \mid x\right] = (E[r \mid x] - E_x[g(x)])^2 + E_x\left[(g(x) - E_x[g(x)])^2\right]$$

*bias*         *variance*

# Estimating Bias and Variance

- $M$ samples $X_i = \{x^t_i, r^t_i\}$, $i = 1, \ldots, M$

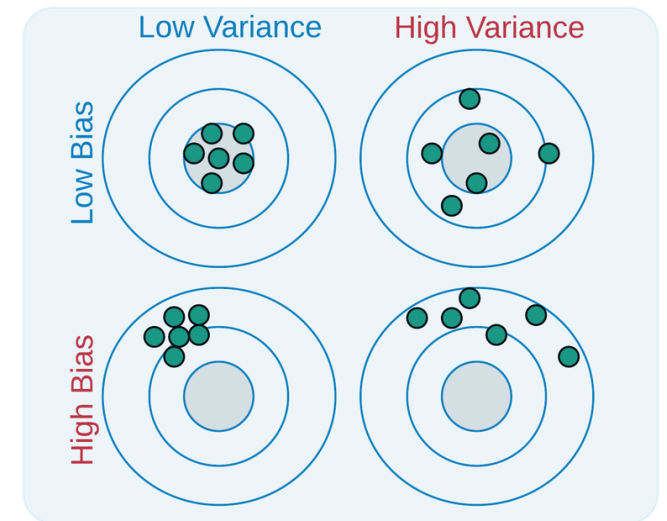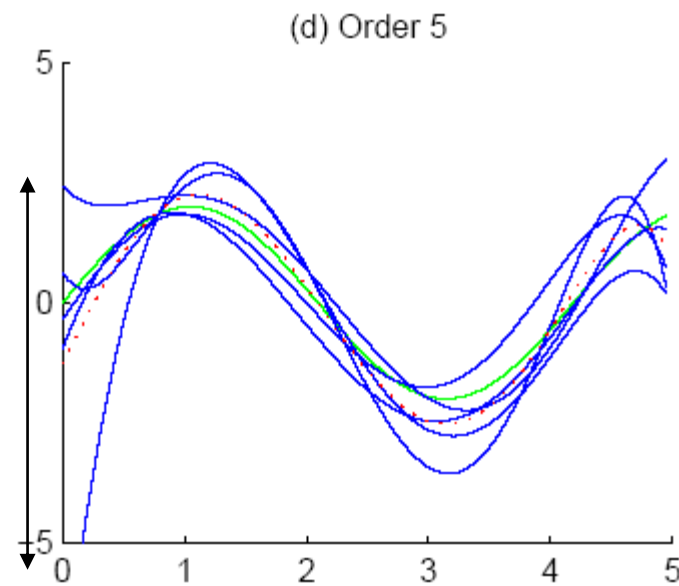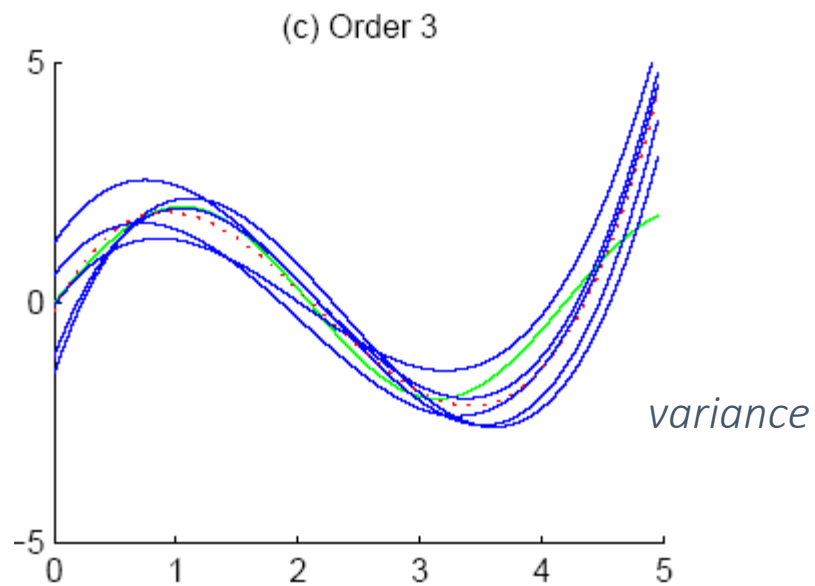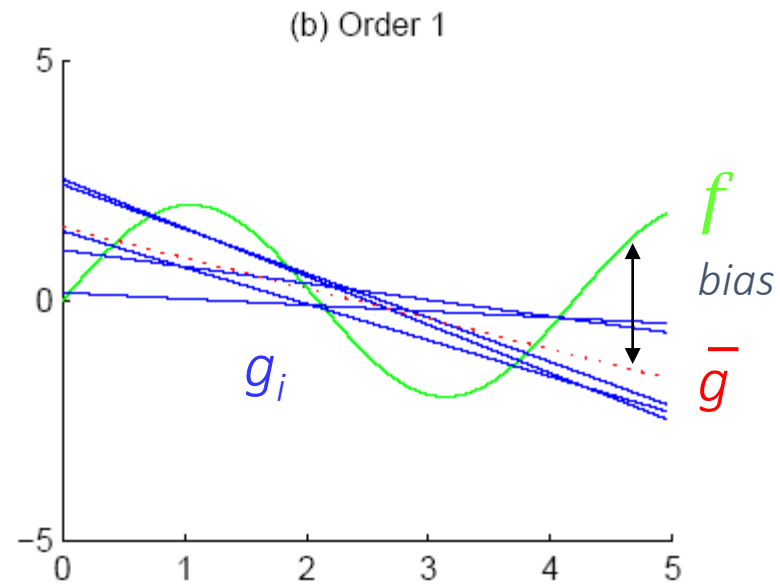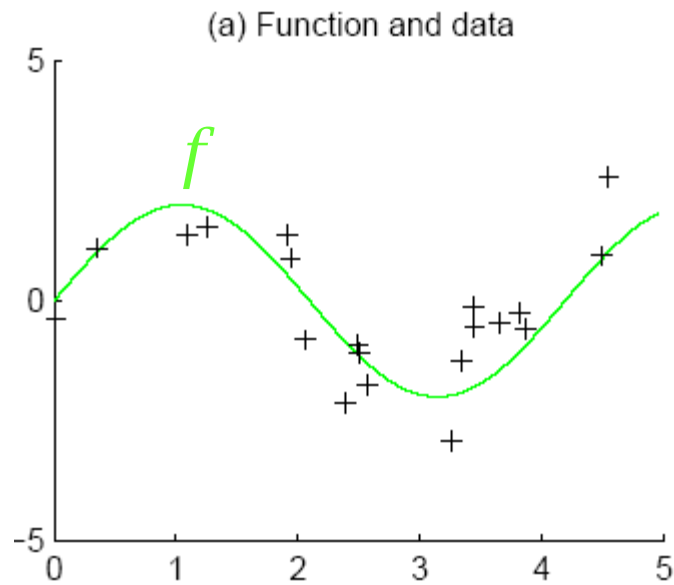  are used to fit $g_i(x)$, $i = 1, \ldots, M$ and $t = 1, \ldots, N$

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

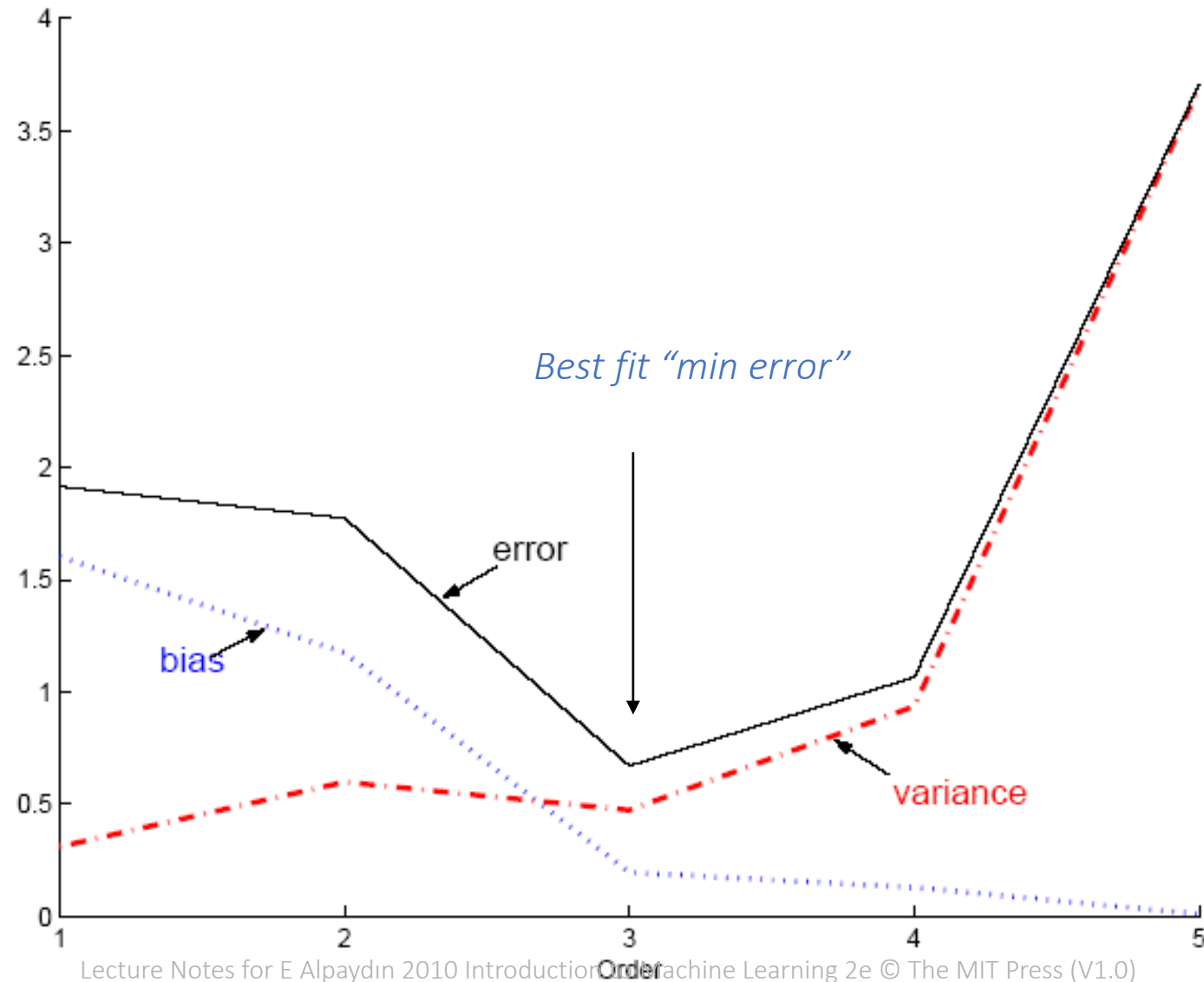$$\text{Variance}(g) = \frac{1}{NM} \sum_t \sum_i \left[ g_i(x^t) - \bar{g}(x^t) \right]^2$$

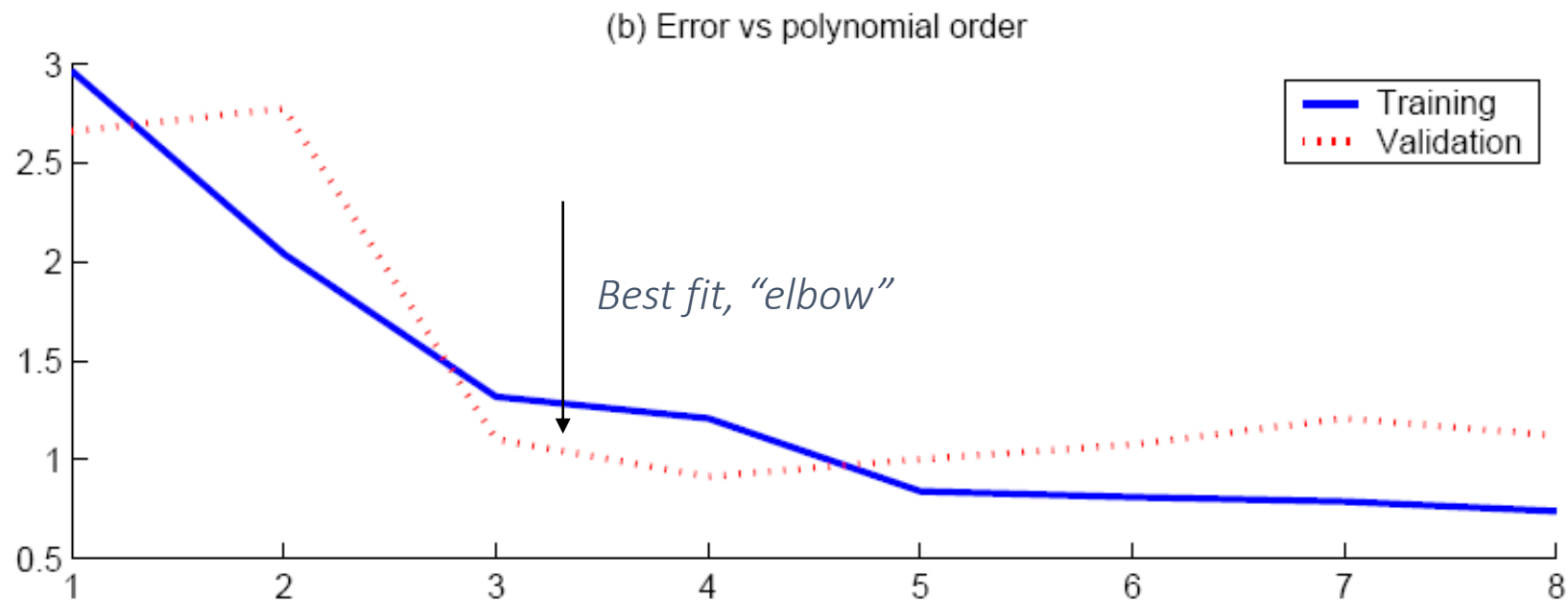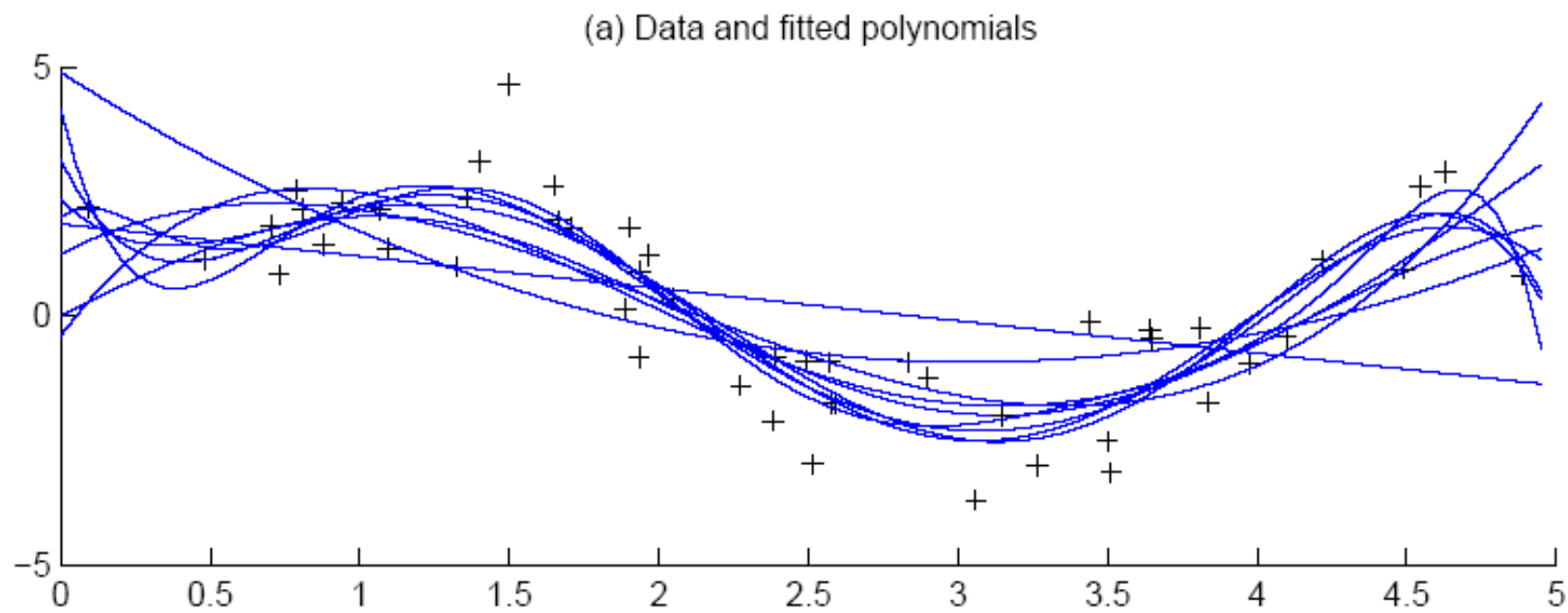$$\bar{g}(x) = \frac{1}{M} \sum_i g_i(x)$$

# Bias/Variance Dilemma

- Example: $g_i(x)=2$ has no variance and high bias

  $g_i(x)= \sum_t r^t_i/N$ 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)
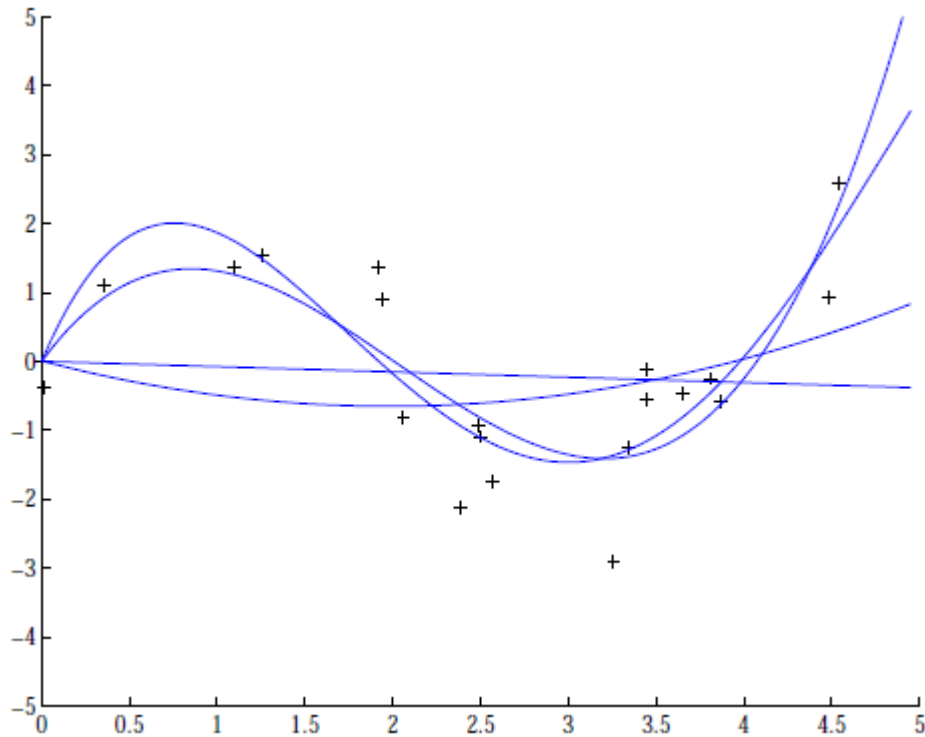
(a) Function and data

(b) Order 1

(c) Order 3

(d) Order 5

Low Variance   High Variance

Low Bias   High Bias

# Polynomial Regression

(a) Data and fitted polynomials

(b) Error vs polynomial order

Best fit, "elbow"

Training
Validation

# 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

$$E(\mathbf{w} \mid \mathcal{X}) = \frac{1}{2} \sum_{t=1}^{N} \left[ y^t - g\left( x^t \mid \mathbf{w} \right) \right]^2 + \lambda \sum_i w_i^2$$

- Ridge Regression

$$R(w) = \|w\|^2 = \sum_i w_i^2$$

- LASSO:

$$R(w) = \|w\|_1 = \sum_i |w_i|$$

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

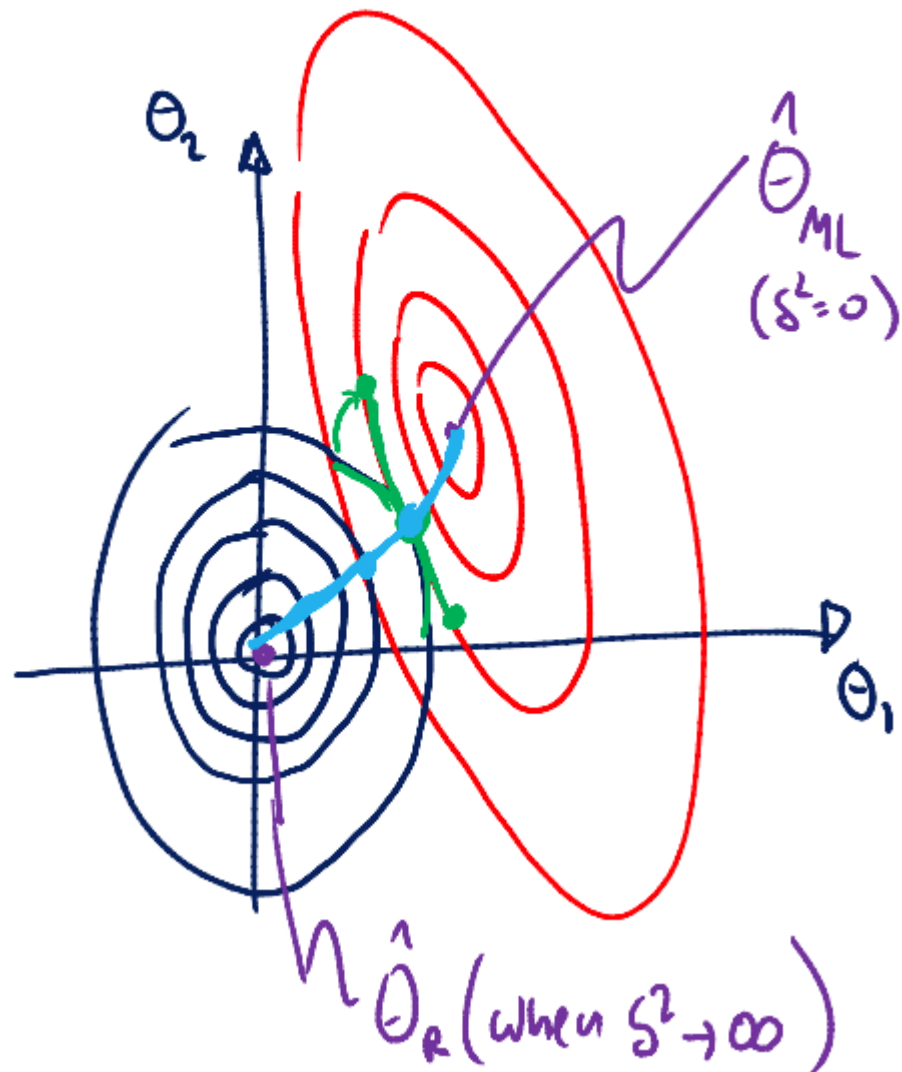- $\nabla \mathcal{L} = -\frac{2}{2} X^T (y - Xw) + \lambda w$

- $-\frac{2}{2} X^T (y - Xw) + \lambda w = 0 \rightarrow X^T y = X^T X w + \lambda w \rightarrow$

- $X^T y = (X^T X + \lambda I) w$

- $\widehat{w} = (X^T X + \lambda I)^{-1} X^T y$

ellipses

$$J(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \delta^2 \boldsymbol{\theta}^T \boldsymbol{\theta}$$



- Image is obtained from Nando Freitas' lecture notes