

# Bias-Var

Slides taken from

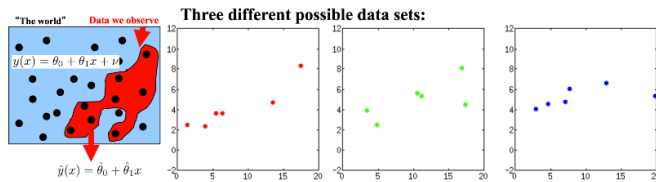
## LEARNING FROM DATA A SHORT COURSE

Yaser S. Abu-Mostafa  
*California Institute of Technology*

Malik Magdon-Ismail  
*Rensselaer Polytechnic Institute*

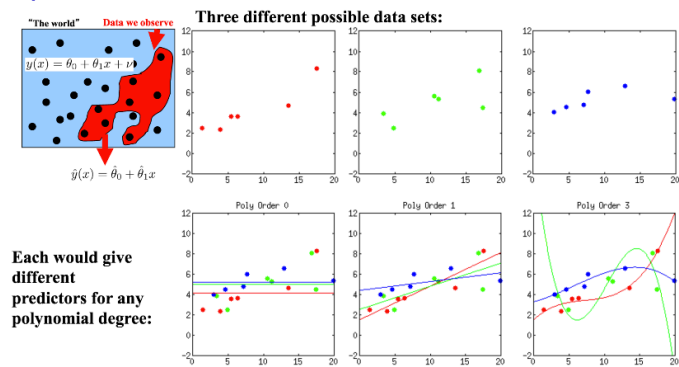
Hsuan-Tien Lin  
*National Taiwan University*

## Bias & variance



Slide by Alexander Ihler

## Bias & variance



Slide by Alexander Ihler

Start with  $E_{\text{out}}$

$$E_{\text{out}}(g^{(\mathcal{D})}) = \mathbb{E}_{\mathbf{x}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right]$$

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} \left[ E_{\text{out}}(g^{(\mathcal{D})}) \right] &= \mathbb{E}_{\mathcal{D}} \left[ \mathbb{E}_{\mathbf{x}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right] \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[ \mathbb{E}_{\mathcal{D}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right] \right] \end{aligned}$$

Now, let us focus on:

$$\mathbb{E}_{\mathcal{D}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right]$$

The average hypothesis

To evaluate  $\mathbb{E}_{\mathcal{D}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right]$

we define the 'average' hypothesis  $\bar{g}(\mathbf{x})$ :

$$\bar{g}(\mathbf{x}) = \mathbb{E}_{\mathcal{D}} \left[ g^{(\mathcal{D})}(\mathbf{x}) \right]$$

Imagine **many** data sets  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$

$$\bar{g}(\mathbf{x}) \approx \frac{1}{K} \sum_{k=1}^K g^{(\mathcal{D}_k)}(\mathbf{x})$$

Using  $\bar{g}(\mathbf{x})$

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right] &= \mathbb{E}_{\mathcal{D}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}) + \bar{g}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right] \\ &= \mathbb{E}_{\mathcal{D}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}) \right)^2 + \left( \bar{g}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right. \\ &\quad \left. + 2 \left( g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}) \right) \left( \bar{g}(\mathbf{x}) - f(\mathbf{x}) \right) \right] \\ &= \mathbb{E}_{\mathcal{D}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}) \right)^2 \right] + \left( \bar{g}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \end{aligned}$$

Bias and variance

$$\mathbb{E}_{\mathcal{D}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right] = \underbrace{\mathbb{E}_{\mathcal{D}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}) \right)^2 \right]}_{\text{var}(\mathbf{x})} + \underbrace{\left( \bar{g}(\mathbf{x}) - f(\mathbf{x}) \right)^2}_{\text{bias}(\mathbf{x})}$$

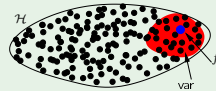
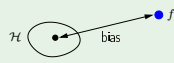
Therefore,  $\mathbb{E}_{\mathcal{D}} \left[ E_{\text{out}}(g^{(\mathcal{D})}) \right] = \mathbb{E}_{\mathbf{x}} \left[ \mathbb{E}_{\mathcal{D}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right] \right]$

$$= \mathbb{E}_{\mathbf{x}} [\text{bias}(\mathbf{x}) + \text{var}(\mathbf{x})]$$

$$= \text{bias} + \text{var}$$

## The tradeoff

$$\text{bias} = \mathbb{E}_{\mathbf{x}} \left[ (\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2 \right] \quad \text{var} = \mathbb{E}_{\mathbf{x}} \left[ \mathbb{E}_{\mathcal{D}} \left[ (g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2 \right] \right]$$


 $\mathcal{H} \uparrow$ 


© 2016 Google. Year 10/10/16 - LFD Lecture 8

9/22

## Example: sine target

$$f: [-1, 1] \rightarrow \mathbb{R} \quad f(x) = \sin(\pi x)$$

Only two training examples!  $N = 2$

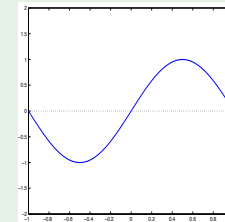
Two models used for learning:

$$\mathcal{H}_0: h(x) = b$$

$$\mathcal{H}_1: h(x) = ax + b$$

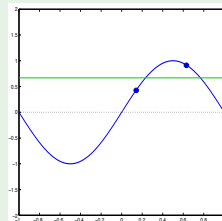
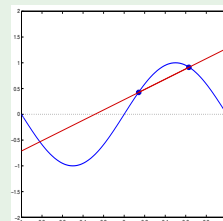
Which is better,  $\mathcal{H}_0$  or  $\mathcal{H}_1$ ?

f



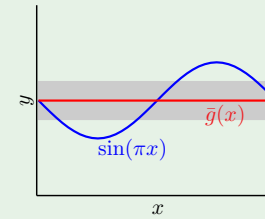
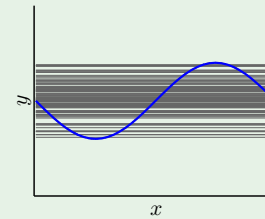
© 2016 Google. Year 10/10/16 - LFD Lecture 8

10/22

Learning -  $\mathcal{H}_0$  versus  $\mathcal{H}_1$  $\mathcal{H}_0$  $\mathcal{H}_1$ 

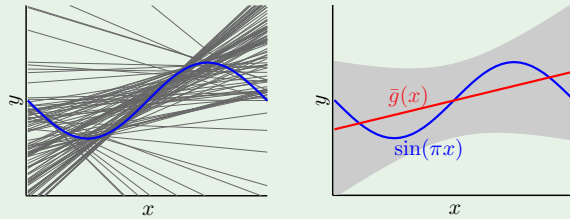
© 2016 Google. Year 10/10/16 - LFD Lecture 8

12/22

Bias and variance -  $\mathcal{H}_0$ 

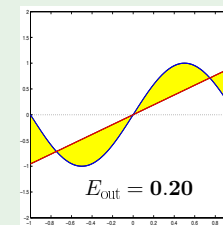
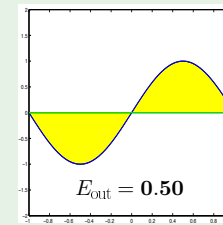
© 2016 Google. Year 10/10/16 - LFD Lecture 8

13/22

Bias and variance -  $\mathcal{H}_1$ 

© Geron: Year Abu-Mostafa - LFD Lecture 8

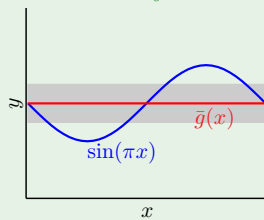
14/22

Approximation -  $\mathcal{H}_0$  versus  $\mathcal{H}_1$  $\mathcal{H}_0$  $\mathcal{H}_1$ 

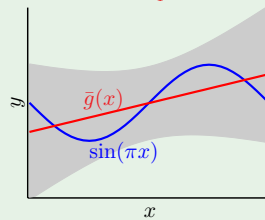
© Geron: Year Abu-Mostafa - LFD Lecture 8

11/22

## and the winner is ...

 $\mathcal{H}_0$  $\mathcal{H}_1$ 

bias = 0.50    var = 0.25



bias = 0.21    var = 1.69

© Geron: Year Abu-Mostafa - LFD Lecture 8

15/22

## Lesson learned

Match the 'model complexity'

to the **data resources**, not to the **target complexity**

© Geron: Year Abu-Mostafa - LFD Lecture 8

16/22

## Outline

- Bias and Variance
- Learning Curves

© 2016 Georgia Tech - Year 10/10/16 - LFD Lecture 8

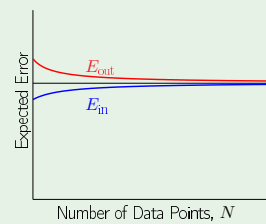
17/22

Expected  $E_{\text{out}}$  and  $E_{\text{in}}$ Data set  $\mathcal{D}$  of size  $N$ Expected out-of-sample error  $\mathbb{E}_{\mathcal{D}}[E_{\text{out}}(g^{(\mathcal{D})})]$ Expected in-sample error  $\mathbb{E}_{\mathcal{D}}[E_{\text{in}}(g^{(\mathcal{D})})]$ How do they vary with  $N$ ?

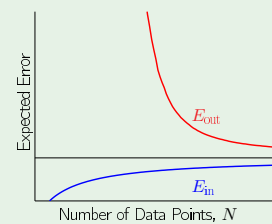
© 2016 Georgia Tech - Year 10/10/16 - LFD Lecture 8

18/22

## The curves



Simple Model

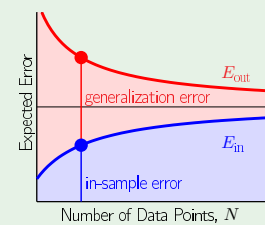


Complex Model

© 2016 Georgia Tech - Year 10/10/16 - LFD Lecture 8

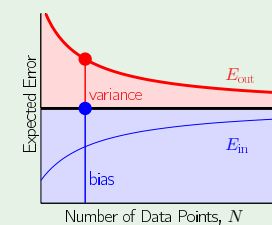
19/22

## VC versus bias-variance



VC analysis

© 2016 Georgia Tech - Year 10/10/16 - LFD Lecture 8



bias-variance

20/22

## Linear regression case

Noisy target  $y = \mathbf{w}^* \mathbf{x} + \text{noise}$

Data set  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$

Linear regression solution:  $\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}$

In-sample error vector =  $X\mathbf{w} - \mathbf{y}$

'Out-of-sample' error vector =  $X\mathbf{w} - \mathbf{y}'$

© 2016 Georgia Tech - Yoram Aloni - LFD Lecture 8

21/22

## Linear regression case

Noisy target  $y = \mathbf{w}^* \mathbf{x} + \text{noise}$

Data set  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$

Linear regression solution:  $\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}$

In-sample error vector =  $X\mathbf{w} - \mathbf{y}$

'Out-of-sample' error vector =  $X\mathbf{w} - \mathbf{y}'$

© 2016 Georgia Tech - Yoram Aloni - LFD Lecture 8

21/22

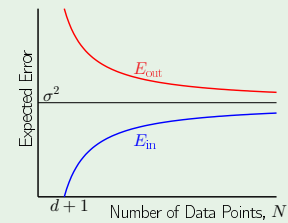
## Learning curves for linear regression

Best approximation error =  $\sigma^2$

Expected in-sample error =  $\sigma^2 \left(1 - \frac{d+1}{N}\right)$

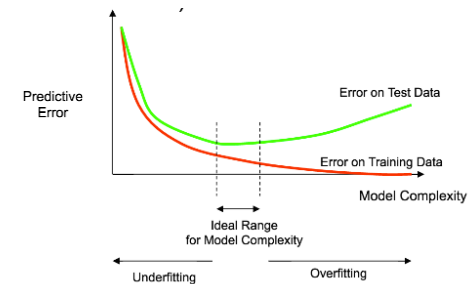
Expected out-of-sample error =  $\sigma^2 \left(1 + \frac{d+1}{N}\right)$

Expected generalization error =  $2\sigma^2 \left(\frac{d+1}{N}\right)$



© 2016 Georgia Tech - Yoram Aloni - LFD Lecture 8

22/22



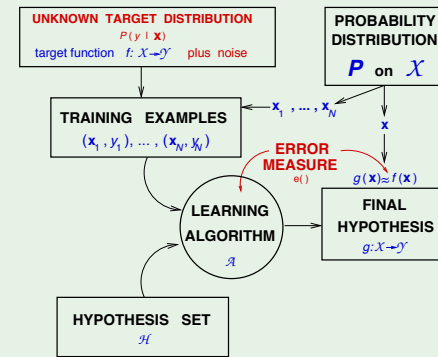
Andrew Ng

### Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its prediction. What should you try next?

- Get more training examples  $\rightarrow$  fixes high variance
- Try smaller sets of features  $\rightarrow$  fixes high variance
- Try getting additional features  $\rightarrow$  fixes high bias
- Try adding polynomial features ( $x_1^2, x_2^2, x_1x_2$ , etc)  $\rightarrow$  fixes high bias.
- Try decreasing  $\lambda \rightarrow$  fixes high bias
- Try increasing  $\lambda \rightarrow$  fixes high variance

### The learning diagram - including noisy target



© Creator: Yasser Abu-Mostafa - LFD Lecture 4