

# ECE595 / STAT598: Machine Learning I

## Lecture 29.1: Bias and Variance - From VC Analysis to Bias-Variance

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering  
Purdue University



# Outline

- Lecture 28 Sample and Model Complexity
- Lecture 29 Bias and Variance
- Lecture 30 Overfit

## Today's Lecture:

- From VC Analysis to Bias-Variance
  - Generalization Bound
  - Bias-Variance Decomposition
  - Interpreting Bias-Variance
- Example
  - 0-th order vs 1-st order model
  - Trade off

# Generalizing the Generalization Bound

## Theorem (Generalization Bound)

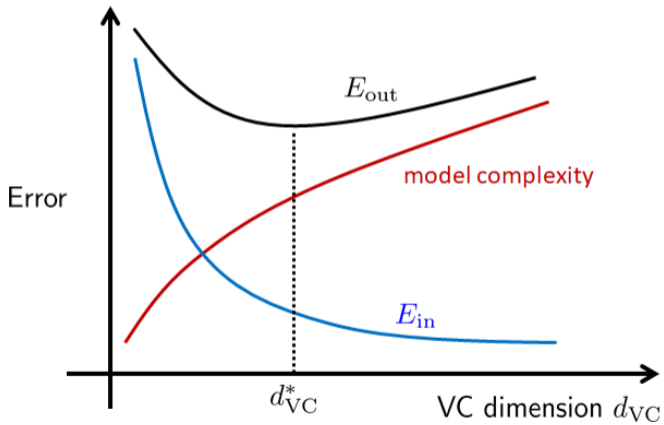
For any tolerance  $\delta > 0$

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{8}{N} \log \frac{4m_{\mathcal{H}}(2N)}{\delta}},$$

with probability at least  $1 - \delta$ .

- $g$ : final hypothesis
- $m_{\mathcal{H}}(N)$ : how complex is your model
- $d_{\text{VC}}$ : parameter defining  $m_{\mathcal{H}}(N) \leq N^{d_{\text{VC}}} + 1$
- Large  $d_{\text{VC}} =$  more complex
- So more difficult to train, and hence require more training samples

# Trade-off Curve



## VC Analysis

- VC analysis is a **decomposition**.
- Decompose  $E_{\text{out}}$  into  $E_{\text{in}}$  and  $\epsilon$ .

$$E_{\text{out}} \leq E_{\text{in}} + \underbrace{\sqrt{\frac{8}{N} \log \frac{4((2N)^{d_{\text{VC}}} + 1)}{\delta}}}_{=\epsilon}$$

- $E_{\text{in}}$  = training error,  $\epsilon$  = penalty of complex model.
- Bias and variance is another decomposition.
- Decompose  $E_{\text{out}}$  into
  - How well can  $\mathcal{H}$  approximate  $f$ ?
  - How well can we zoom in a good  $h$  in  $\mathcal{H}$ ?
- Roughly speaking we will have

$$E_{\text{out}} = \text{bias} + \text{variance}$$

## From VC Analysis to Bias-Variance

- In **VC analysis** we define the out-sample error as

$$E_{\text{out}}(g) = \mathbb{P}[g(\mathbf{x}) \neq f(\mathbf{x})]$$

- Let  $B = \{g(\mathbf{x}) \neq f(\mathbf{x})\}$  be the bad event.  $B \in \{0, 1\}$ .
- Then this is equal to

$$\begin{aligned} E_{\text{out}}(g) &= \mathbb{P}[B = 1] \\ &= 1 \cdot \mathbb{P}[B = 1] + 0 \cdot \mathbb{P}[B = 0] \\ &= \mathbb{E}[B]. \end{aligned}$$

- So  $E_{\text{out}}(g)$  can be written as

$$E_{\text{out}}(g) = \mathbb{E}_{\mathbf{x}}[\mathbf{1}\{g(\mathbf{x}) \neq f(\mathbf{x})\}].$$

- Expectation taken over all  $\mathbf{x} \sim p(\mathbf{x})$ .

## Changing the Error Measure

- In **VC analysis** we define the out-sample error as

$$E_{\text{out}}(g) = \mathbb{E}_{\mathbf{x}} \left[ \mathbf{1}\{g(\mathbf{x}) \neq f(\mathbf{x})\} \right]$$

- Expectation of a **0-1 loss**.
- In **Bias-variance** analysis we define the out-sample error as

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

- Expectation of a **square loss**.
- Square loss is differentiable.

## Dependency on Training Set

- In VC analysis we define the out-sample error as

$$E_{\text{out}}(g^{(\mathcal{D})}) = \mathbb{E}_{\mathbf{x}} \left[ \mathbf{1}\{g^{(\mathcal{D})}(\mathbf{x}) \neq f(\mathbf{x})\} \right]$$

- The final hypothesis depends on  $\mathcal{D}$ .
- If you use a different  $\mathcal{D}$ , your  $g$  will be different.
- In Bias-variance analysis we define the out-sample error as

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

- Why did we skip  $\mathcal{D}$  in VC analysis?
  - Hoeffding bound is uniform for **all**  $\mathcal{D}$
  - So it does not matter which  $\mathcal{D}$  you used to generate  $g$
  - Not true for bias-variance



## Averaging over all $\mathcal{D}$

- To account for all the possible  $\mathcal{D}$ 's, compute the expectation and define the expected out-sample error.

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

- $E_{\text{out}}(g^{(\mathcal{D})})$ : Out-sample error for the particular  $g$  found from  $\mathcal{D}$
- $\mathbb{E}_{\mathcal{D}} [E_{\text{out}}(g^{(\mathcal{D})})]$ : Out-sample error averaged over all possible  $\mathcal{D}$ 's
- VC trade-off is a “worst case” analysis
  - Uniform bound on every  $\mathcal{D}$
- Bias-variance trade-off is an “average” analysis
  - Average over different  $\mathcal{D}$ 's

## Decomposing $\mathbb{E}_{\text{out}}(g^{(\mathcal{D})})$

- To account for all the possible  $\mathcal{D}$ 's, compute the expectation and define the expected out-sample error.

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

- Let us do some calculation

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

## The Average $\bar{g}(\mathbf{x})$

- The decomposition gives

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

- We define the term

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

- The asymptotic limit of the estimate

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

- $g^{(\mathcal{D}_k)}$  are inside the hypothesis set. But  $\bar{g}$  is *not* necessarily inside.

## Bias and Variance

- Do some additional calculation

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

- Define two terms

$$\text{bias}(\mathbf{x}) \stackrel{\text{def}}{=} (\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2,$$

$$\text{var}(\mathbf{x}) \stackrel{\text{def}}{=} \mathbb{E}_{\mathcal{D}}[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2].$$

## Bias and Variance

- The decomposition:

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

- Define two terms

$$\begin{aligned} \text{bias}(\mathbf{x}) & \stackrel{\text{def}}{=} (\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2, \\ \text{var}(\mathbf{x}) & \stackrel{\text{def}}{=} \mathbb{E}_{\mathcal{D}}[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2]. \end{aligned}$$

- Take expectation

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

# Bias and Variance Decomposition

- The decomposition:

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

- This gives

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} \left[ \mathbb{E}_{\text{out}}(g^{(\mathcal{D})}) \right] &= \mathbb{E}_{\mathbf{x}}[\text{bias}(\mathbf{x}) + \text{var}(\mathbf{x})] \\ &= \text{bias} + \text{var} \end{aligned}$$

## Interpreting the Bias-Variance

- The decomposition:

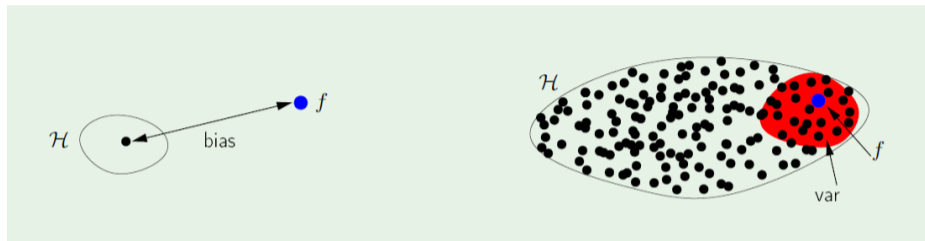
$$\begin{aligned} & \mathbb{E}_{\mathcal{D}} \left[ \mathbb{E}_{\text{out}}(g^{(\mathcal{D})}) \right] \\ = & \mathbb{E}_{\mathbf{x}} \left[ \underbrace{\mathbb{E}_{\mathcal{D}} \left[ g^{(\mathcal{D})}(\mathbf{x})^2 \right]}_{\mathbb{E}_{\mathcal{D}}[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2]}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

- The two terms:

$$\begin{aligned} \text{bias}(\mathbf{x}) & \stackrel{\text{def}}{=} (\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2, \\ \text{var}(\mathbf{x}) & \stackrel{\text{def}}{=} \mathbb{E}_{\mathcal{D}}[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2]. \end{aligned}$$

- $\text{bias}(\mathbf{x})$ : How close is the **average function**  $\bar{g}$  to the target
- $\text{var}(\mathbf{x})$ : How much **uncertainty** you have around  $\bar{g}$

# Model Complexity



- The bias and variance are

$$\text{bias}(\mathbf{x}) \stackrel{\text{def}}{=} (\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2,$$
$$\text{var}(\mathbf{x}) \stackrel{\text{def}}{=} \mathbb{E}_{\mathcal{D}}[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2].$$

- If you have a simple  $\mathcal{H}$ , then large bias but small variance
- If you have a complex  $\mathcal{H}$ , then small bias but large variance