

Machine Learning from Data

Lecture 7: Spring 2021

Today's Lecture

- Approximation Vs Generalization
 - VC Dimension
 - Bias-Variance Analysis
 - Learning Curve

Theory of Learning and its relevance.

- We created a link between E_{out} and E_{in}

$$\mathbb{P} [|\boldsymbol{E}_{\text{in}}(\boldsymbol{g}) - \boldsymbol{E}_{\text{out}}(\boldsymbol{g})| > \epsilon] \leq 4m_{\mathcal{H}}(\mathbf{2N})e^{-\epsilon^2N/8}, \qquad \text{for any } \epsilon > 0.$$

$$\mathbb{P}[|E_{\text{in}}(g)-E_{\text{out}}(g)|>\epsilon]\leq 2|\mathcal{H}|e^{-2\epsilon^2N} \quad \leftarrow \text{finite } \mathcal{H}$$

$$\mathbb{P} [|\boldsymbol{E}_{\text{in}}(\boldsymbol{g}) - \boldsymbol{E}_{\text{out}}(\boldsymbol{g})| \leq \epsilon] \geq 1 - 4m_{\mathcal{H}}(\mathbf{2N})e^{-\epsilon^2N/8}, \qquad \text{for any } \epsilon > 0.$$

$$\mathbb{P}[|E_{\text{in}}(g)-E_{\text{out}}(g)|\leq\epsilon]\geq 1-2|\mathcal{H}|e^{-2\epsilon^2N} \quad \leftarrow \text{finite } \mathcal{H}$$

$$\boldsymbol{E}_{\text{out}}(\boldsymbol{g}) \leq \boldsymbol{E}_{\text{in}}(\boldsymbol{g}) + \sqrt{\frac{8}{N} \log \frac{4m_{\mathcal{H}}(\mathbf{2N})}{\delta}}, \qquad \text{w.p. at least } \mathbf{1} - \delta.$$

$$E_{\text{out}}(g)\leq E_{\text{in}}(g)+\sqrt{\frac{1}{2N} \log \frac{2|\mathcal{H}|}{\delta}} \quad \leftarrow \text{finite } \mathcal{H}$$

$$m_{\mathcal{H}}(N) \leq \sum_{i=1}^{\boldsymbol{k}-1} \binom{N}{i} \leq N^{\boldsymbol{k}-1} + 1 \qquad \boldsymbol{k} \text{ is a break point.}$$

The VC Dimension

-

Summarize

$$m_{\mathcal{H}}(N) \sim N^{k-1}$$

The tightest bound is obtained with the smallest break point k^* .

Definition [VC Dimension] $d_{\text{VC}} = k^* - 1$.

The VC dimension is the largest N which can be shattered ($m_{\mathcal{H}}(N) = 2^N$).

$N \leq d_{\text{VC}}$: \mathcal{H} could shatter your data (\mathcal{H} can shatter some N points).

$N > d_{\text{VC}}$: N is a break point for \mathcal{H} ; \mathcal{H} cannot possibly shatter your data.

$$m_{\mathcal{H}}(N) \leq N^{d_{\text{VC}}} + 1 \sim N^{d_{\text{VC}}}$$

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + O\left(\sqrt{\frac{d_{\text{VC}} \log N}{N}}\right)$$

Examples

-

Convex Sets

-

Summarize

	N						#Param	d_{VC}
	1	2	3	4	5	...		
2-D perceptron	2	4	8	14	...		3	3
1-D pos. ray	2	3	4	5	...		1	1
2-D pos. rectangles	2	4	8	16	$< 2^5$...	4	4
pos. convex sets	2	4	8	16	32	...	∞	∞

There are models with few parameters but infinite d_{VC} .

There are models with redundant parameters but small d_{VC} .



Sample Complexity?

Set the error bar at ϵ .

$$\epsilon = \sqrt{\frac{8}{N} \ln \frac{4((2N)^{d_{\text{vc}}} + 1)}{\delta}}$$

Solve for N :

$$N = \frac{8}{\epsilon^2} \ln \frac{4((2N)^{d_{\text{vc}}} + 1)}{\delta} = O(d_{\text{vc}} \ln N)$$

Example. $d_{\text{vc}} = 3$; error bar $\epsilon = 0.1$; confidence 90% ($\delta = 0.1$).
A simple iterative method works well. Trying $N = 1000$ we get

$$N \approx \frac{1}{0.1^2} \log \left(\frac{4(2000)^3 + 4}{0.1} \right) \approx 21192.$$

We continue iteratively, and converge to $N \approx 30000$.

If $d_{\text{vc}} = 4$, $N \approx 40000$; for $d_{\text{vc}} = 5$, $N \approx 50000$.

($N \propto d_{\text{vc}}$, but gross overestimates)

Practical Rule of Thumb: $N = 10 \times d_{\text{vc}}$

Theory Vs Practice

-

VC Bound Quantifies Approximation Vs. Generalization

$d_{\text{VC}} \uparrow \implies$ better chance of **approximating** f ($E_{\text{in}} \approx 0$).

$d_{\text{VC}} \downarrow \implies$ better chance of **generalizing** to out of sample ($E_{\text{in}} \approx E_{\text{out}}$).

$$E_{\text{out}} \leq E_{\text{in}} + \Omega(d_{\text{VC}}).$$

Bias Variance Trade-off

1. How well *can* the learning approximate f .

... as opposed to how well *did* the learning approximate f in-sample (E_{in}).

2. How close can you get to that approximation with a finite data set.

... as opposed to how close is E_{in} to E_{out} .

Bias-variance analysis applies to squared errors (classification and regression)

Bias-variance analysis can take into account the *learning algorithm*

Different learning algorithms can have different E_{out} when applied to the same \mathcal{H} !

A simple learning problem

-



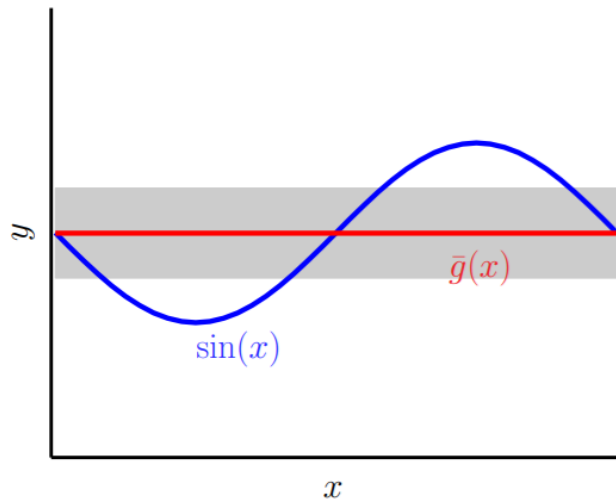
Problem continued...

-

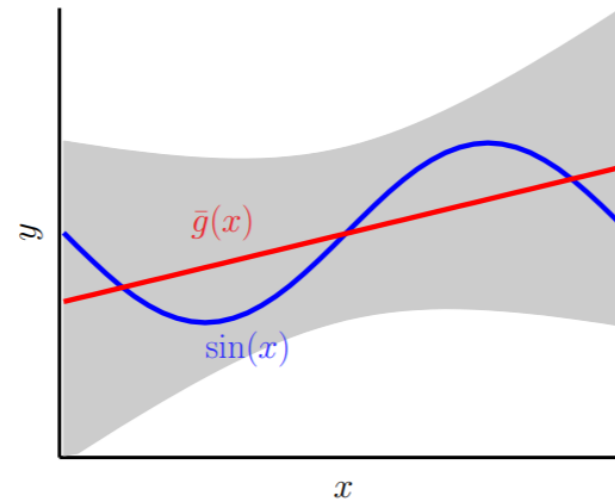
Expected Behavior with Datasets

-

Which one is better?



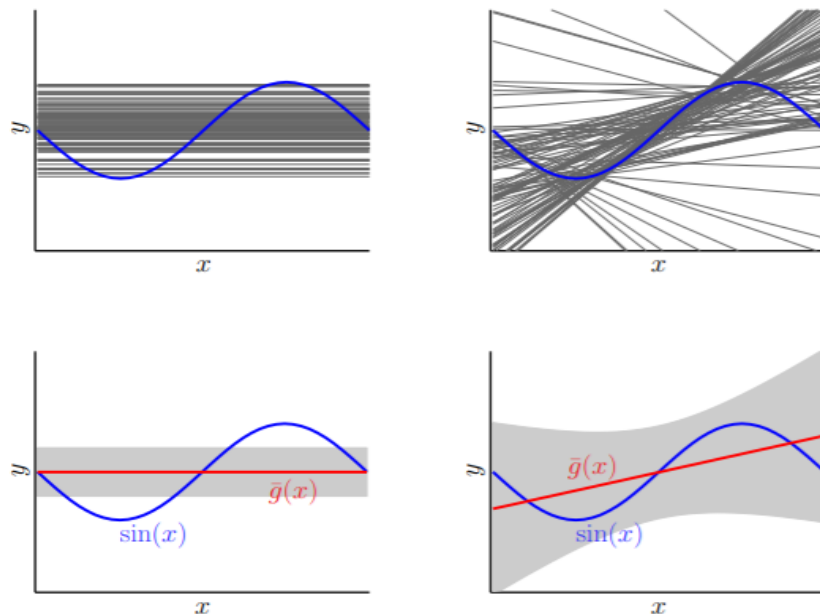
$$\begin{array}{l} \mathcal{H}_0 \\ \text{bias} = 0.50 \\ \text{var} = 0.25 \\ \hline E_{\text{out}} = 0.75 \quad \checkmark \end{array}$$



$$\begin{array}{l} \mathcal{H}_1 \\ \text{bias} = 0.21 \\ \text{var} = 1.69 \\ \hline E_{\text{out}} = 1.90 \end{array}$$

Data

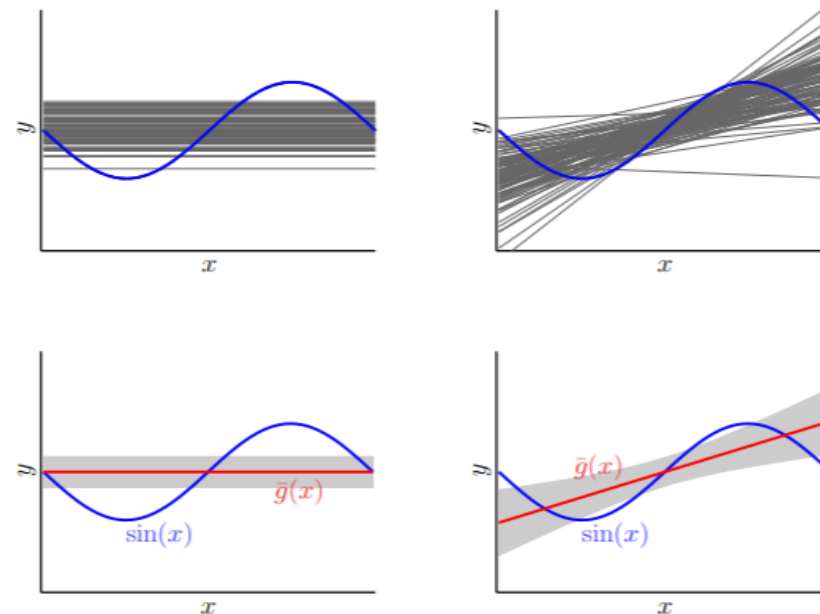
2 Data Points



$$\begin{aligned} &\mathcal{H}_0 \\ &\text{bias} = 0.50; \\ &\text{var} = 0.25. \\ \hline &E_{\text{out}} = 0.75 \quad \checkmark \end{aligned}$$

$$\begin{aligned} &\mathcal{H}_1 \\ &\text{bias} = 0.21; \\ &\text{var} = 1.69. \\ \hline &E_{\text{out}} = 1.90 \end{aligned}$$

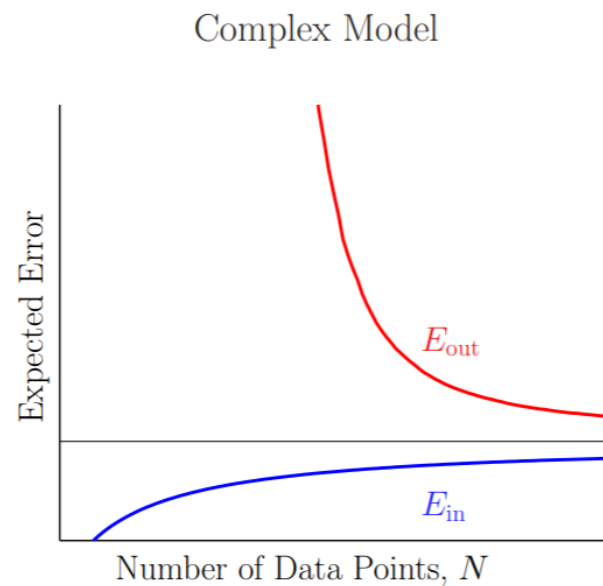
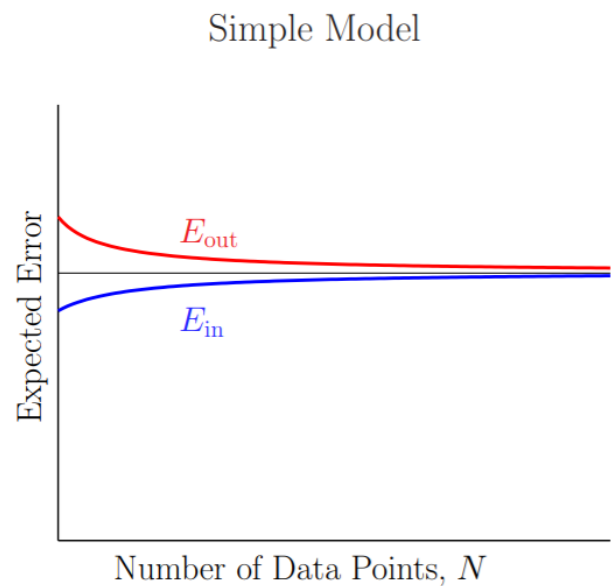
5 Data Points



$$\begin{aligned} &\mathcal{H}_0 \\ &\text{bias} = 0.50; \\ &\text{var} = 0.1. \\ \hline &E_{\text{out}} = 0.6 \end{aligned}$$

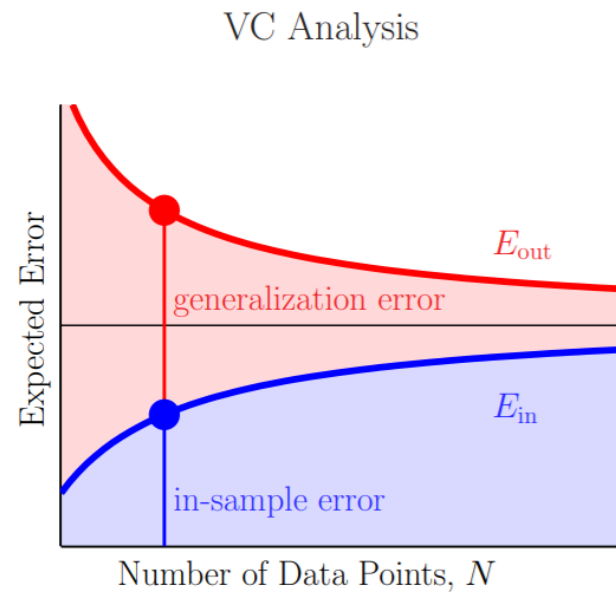
$$\begin{aligned} &\mathcal{H}_1 \\ &\text{bias} = 0.21; \\ &\text{var} = 0.21. \\ \hline &E_{\text{out}} = 0.42 \quad \checkmark \end{aligned}$$

Learning Curve

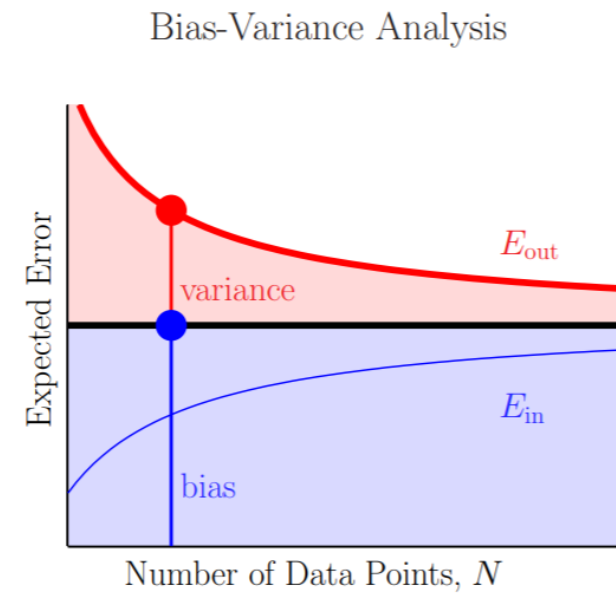


$$E_{\text{out}} = \mathbb{E}_{\mathbf{x}} [E_{\text{out}}(\mathbf{x})]$$

Comparison



Pick \mathcal{H} that can generalize and has a good chance to fit the data



Pick $(\mathcal{H}, \mathcal{A})$ to approximate f and not behave wildly after seeing the data

Thanks!