

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}

Universel Sample
Data $\left(\begin{array}{c} \text{Data} \\ \text{Unknown} \end{array} \right)$ $\left(\begin{array}{c} E_{in} \\ \text{see} \end{array} \right)$ E_{out}

$$E_{out} \leq E_{in} + \sqrt{\frac{\delta}{N} \ln 4(m_h(2x))}$$

(growth function) $\rightarrow \epsilon$

i) fail or succeed with high probability.

2-step learning process:

- i) $E_{in} \approx 0$ and $E_{out} \approx E_{in}$ ✓ ideal
- ii) $E_{in} \gg 0$ and $E_{out} \approx E_{in}$ Not good
- iii) $E_{in} \approx 0$ $E_{out} \gg 0$ Theory helps us avoid this

i) Time tested
ii) link between E_{in} & E_{out}
iii) Relevance in practice.

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

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

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

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

$n, f, P(x), A$

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

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

w.p. at least $1 - \delta$

$$m_{\mathcal{H}}(N) \leq \sum_{i=1}^{k-1} \binom{N}{i} \leq N^{k-1} + 1 \quad (\text{Polynomial})$$

k is a break point.

2 kinds of N sets:

1) Bad $m_k(N) = 2^N \quad \forall N \geq 1$, $\mathcal{R} = \sqrt{\frac{8 \log f(2^N)}{N}}$

$E_{\text{in}} \neq E_{\text{out}}$

$\nexists 0$ expand

2) $m_n(N) \rightarrow$ polynomially

$m_n(N)$ $\leq \sum_{i=0}^{k-1} \binom{N}{i} \leq N^{k-1} + 1 \simeq N^{k-1}$

$\mathcal{R} = \sqrt{\frac{8 \log f(2N)^{k-1}}{N}}$

$\mathcal{R} = \sqrt{\frac{8(k-1) \log(n)}{N}} \rightarrow 0$

The VC Dimension

- The VC dimension of a \mathcal{H} is one less than the cheapest breakpoints (smallest)

$$d_{VC} = k^* - 1$$

$d_{VC} \rightarrow$ max. no. of data points that can be shattered.

If $N \leq d_{VC} \rightarrow \mathcal{H}$ can shatter data

If $N > d_{VC} \rightarrow \mathcal{H}$ cannot shatter your data

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.

$$O(\sqrt{d_{\text{vc}} \log N})$$

grows

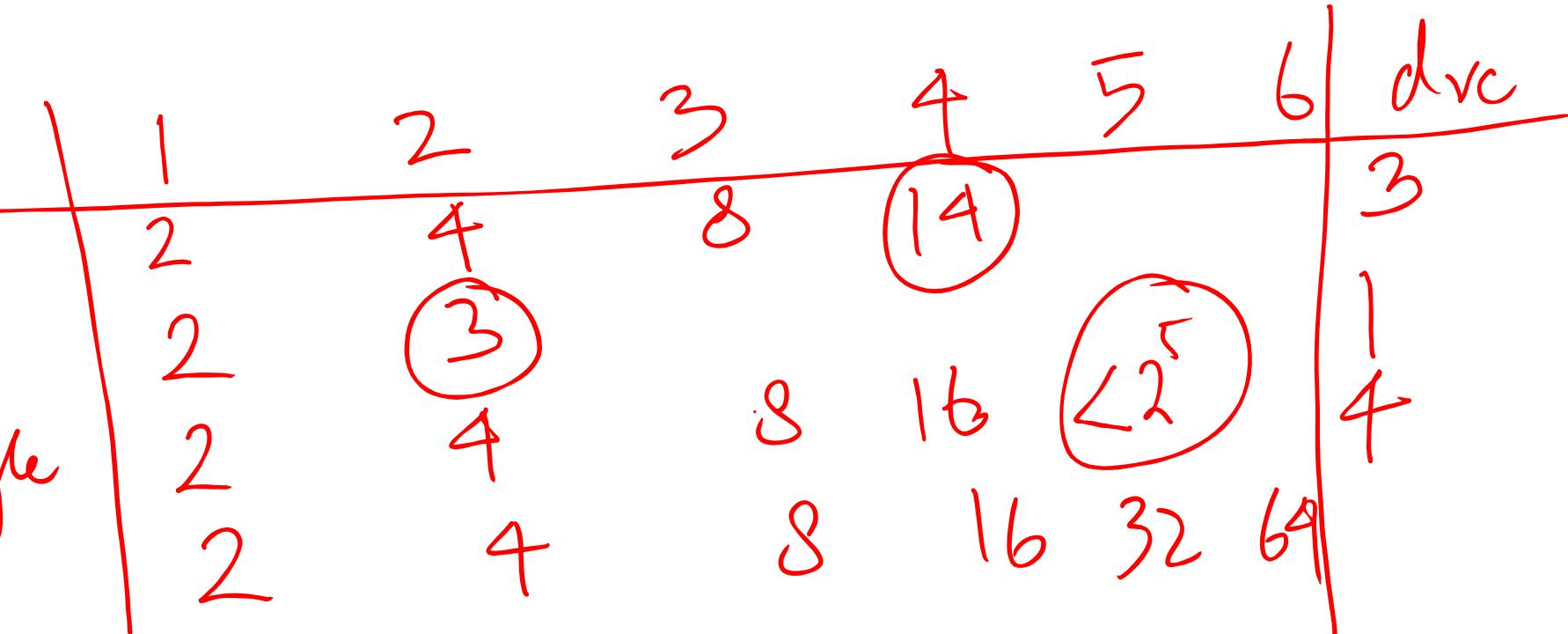
$$m_{\mathcal{H}}(N) \leq N^{d_{\text{vc}}} + 1 \sim N^{d_{\text{vc}}}$$

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

SV
(for good
 N)

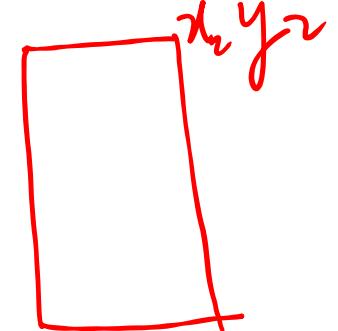
Examples

- $\frac{N}{2}$ 2-d perception
- 1-d tree ray
- 2-d tree rectangle
- 2^n



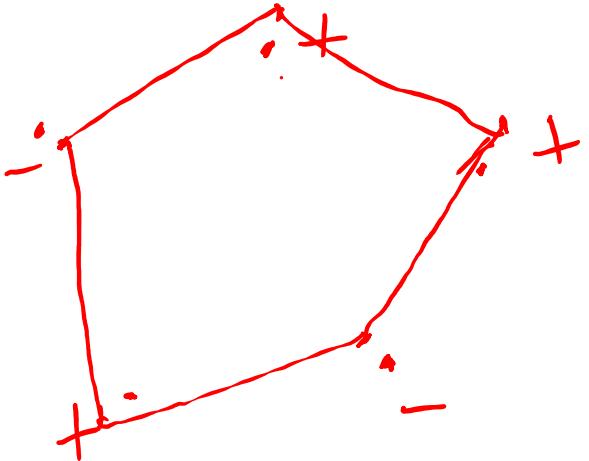
2-d perceptions	w_0, w_1, w_2	3
1-d tree ray	w_0 (threshold)	1
2-d tree rectangle	x_1, y_1, x_2, y_2	4

* Relationship we cannot always guarantee x_1, y_1

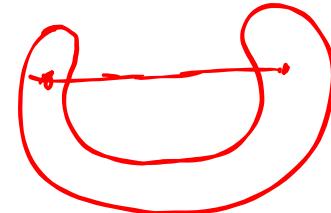
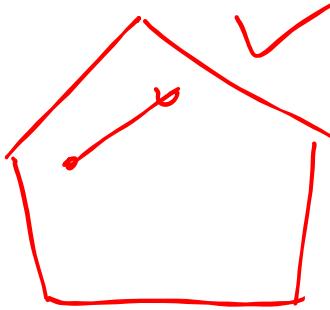


Convex Sets

- the convex sets.



(n)



$d_{VC} = \mathcal{D}$
Not useful for learning.

Summarize

*Not useful
in practice*

	1	2	3	4	5	...	#Param	d_{VC}
	N							
✓ 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} .

Perception : $h(\bar{x}) = \text{sign}(\bar{w}^T \bar{x})$

d -dimensional we have w_0, w_1, \dots, w_d
 $= d+1$ parameters.

Theorem : The VC-dimension of the perception

$$d_{VC}(\text{perception}) = d+1$$

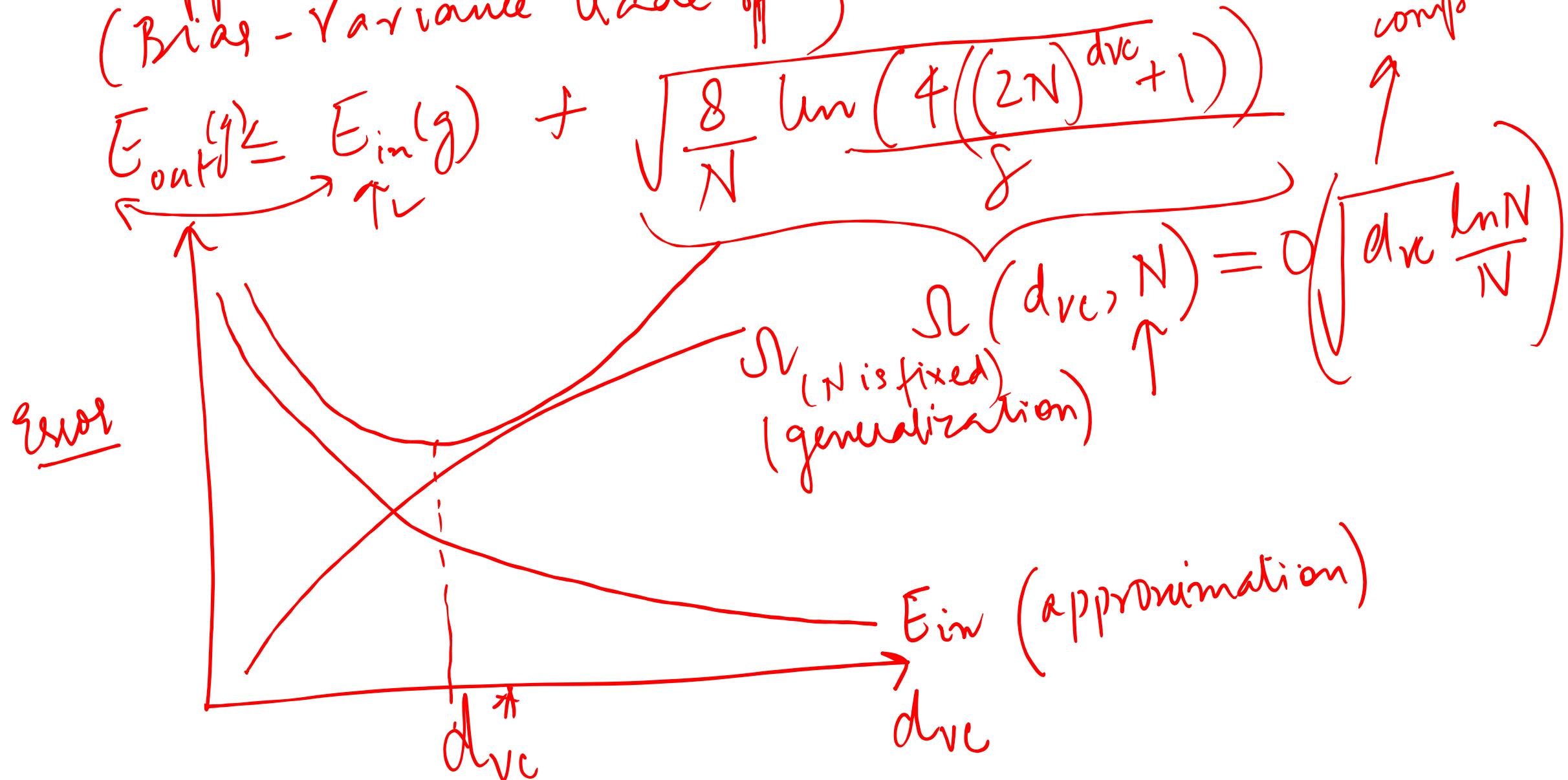
$$\bar{E}_{out} \leq \bar{E}_{in} + O\left(\sqrt{\frac{d_{VC} \ln N}{N}}\right)$$

Proof : (Ex: 2.4 - HW problem)

- i) Show that $d_{VC} \geq d+1 \rightarrow$ Any $(d+1)$ points $\bar{d}_{VC} = \max.$
- ii) Show that $d_{VC} \leq d+1$
 - \rightarrow $d+2$ points.
 - \hookrightarrow linear algebra
(hint)

no. of
data points
you can
separate

Approximation VS Generalization (Bias - Variance Trade-off)



H (neural network) $\rightarrow d_{VC} \rightarrow$ picking some

SAMPLE COMPLEXITY

$\epsilon \rightarrow$ generalization error

$\delta \rightarrow$ confidence parameter.

$$\epsilon = \sqrt{\frac{8}{N} \ln \frac{4(2N)}{\delta}}^{d_{VC}}$$

$$N$$

$E_{in} \simeq E_{out}$

Sample Complexity?

Set the error bar at ϵ .

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

Solve for N :

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

$N = 10 \times d_{VC}$ (simpler)

$N = 100 \times d_{VC}$

$N = 1000 \times d_{VC}$ greater

z-d exception $\downarrow 8$

Example. $d_{VC} = 3$; 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 $d_{VC} = 4$, $N \approx 40000$; for $d_{VC} = 5$, $N \approx 50000$.

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

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

Theory Vs Practice

• Practice tells us $N \downarrow$

• \checkmark VC

general
 M -sets, LA,
Input (X), P^d ,
 f

VC, $N \approx 10,000$

i) Hoeffding Inequality (M -sch)
Accomodate for general

ii) $m_n(N)$
Any \uparrow worst case
dataset

iii) $m_n(N) \leq [B(N, k)] \leq N^{d_{VC} + 1}$
 $d_{VC} + 1$ \uparrow technical slack

VC Bound Quantifies Approximation Vs. Generalization

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

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

i) N -fixed $\xrightarrow{d_{VC}} N$ -size.

ii) N -size $\xrightarrow[\text{fixed}]{d_{VC}} N$ -choose

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

model complexity.

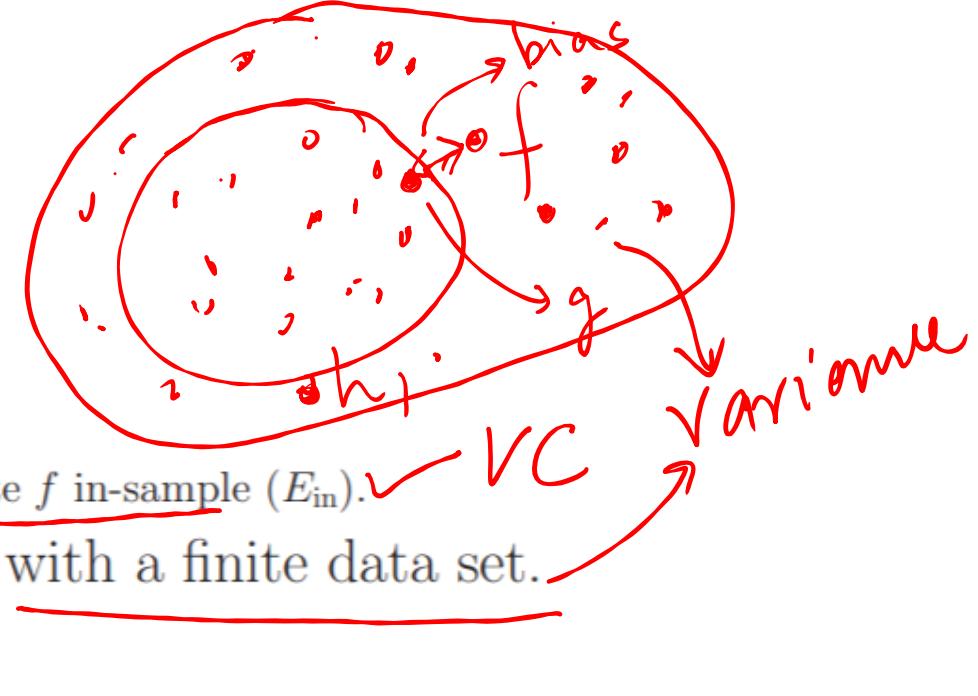
Bias Variance Trade-off

1. How well can the learning approximate f .
Bias

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

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} . *VC*



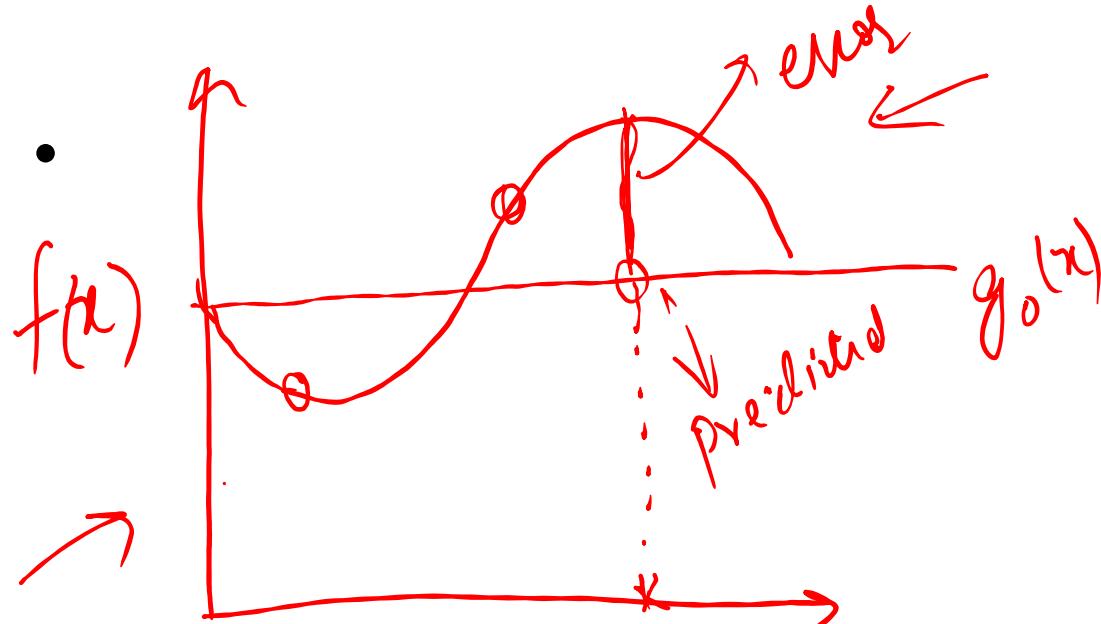
regression problems

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

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

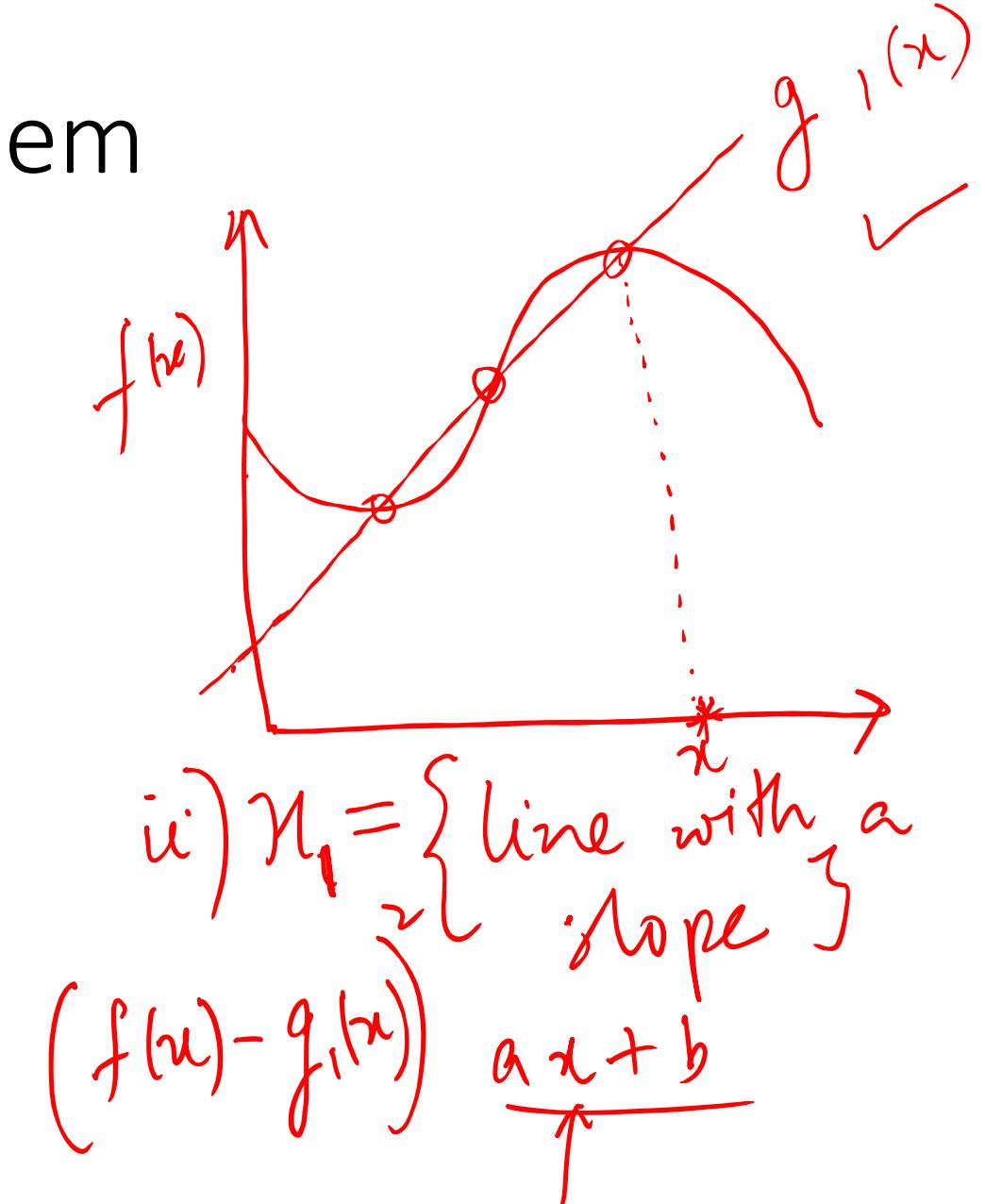
Different learning algorithms can have different E_{out} when applied to the same \mathcal{H} !

A simple learning problem

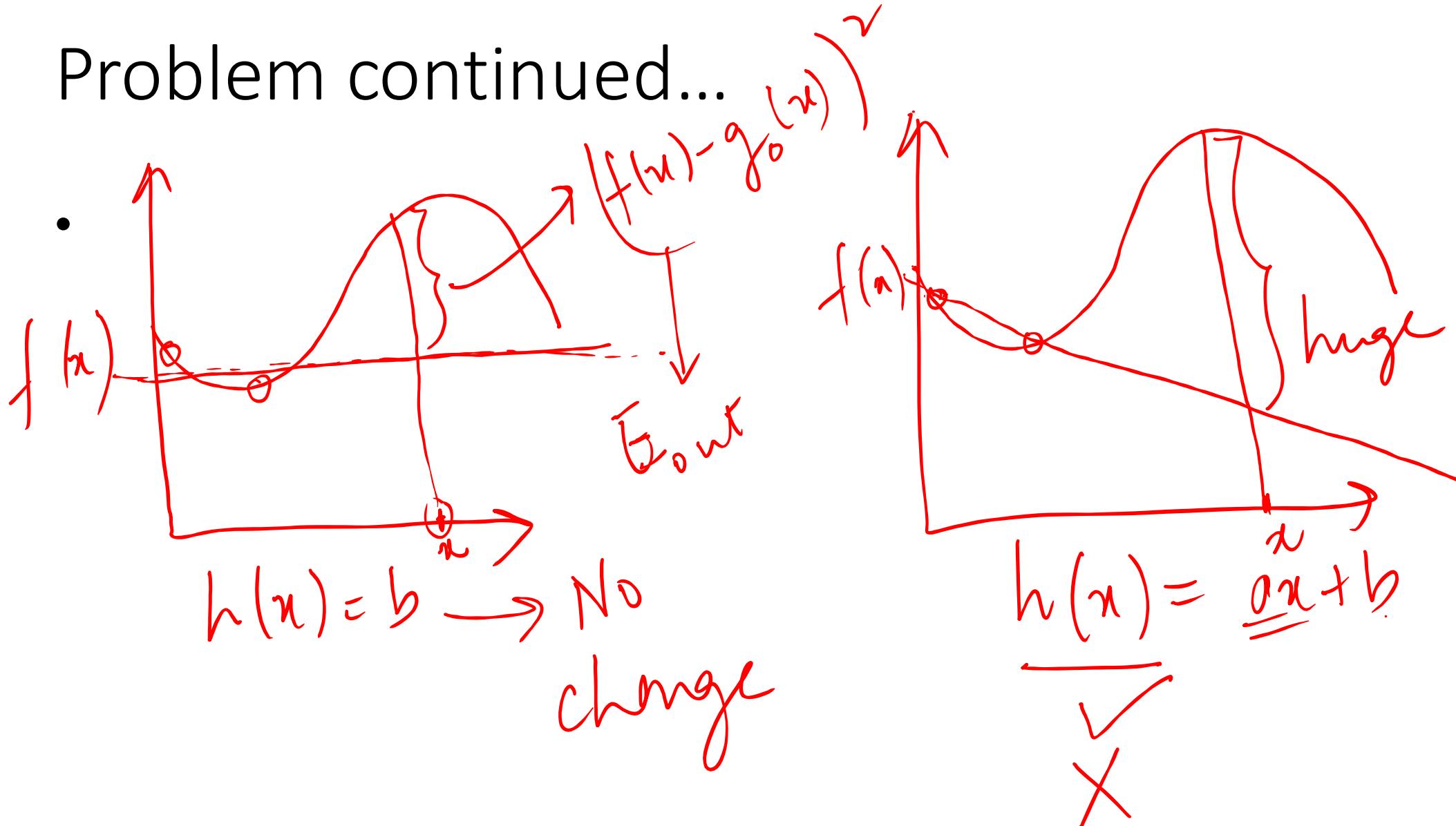


i) $\mathcal{H}_0 = \{\text{flat line}\}$ $h(x) = b$

$$(f(x) - g_0(x))^2$$



Problem continued...



Expected Behavior with Datasets

- $D_1 \rightarrow g^{D_1} \rightarrow g^{D_1}(x)$
- $D_2 \rightarrow g^{D_2} \rightarrow g^{D_2}(x)$
- $D_3 \rightarrow g^{D_3}$
- \vdots
- $D_M \rightarrow g^{D_M} \rightarrow g^{D_M}(x)$

$$\bar{g}(x) = \frac{1}{M} \sum_{i=1}^M g^{D_i}(x)$$

$$\left. \begin{array}{l} \rightarrow E_D [g^D(x)] \rightarrow \text{expected prediction} \\ \text{Var}(x) = \frac{1}{M} \sum_{i=1}^M (g^{D_i}(x) - \bar{g}(x))^2 \\ E_D [(g^D(x) - \bar{g}(x))^2] \\ \downarrow \\ \text{Variance of prediction} \end{array} \right\}$$

$$(f(x) - g^{D_i}(x))^2 \rightarrow \text{Avg. error} = \frac{1}{M} \sum_{i=1}^M (f(x) - g^{D_i}(x))^2$$

Result from probability

↑
constant
steps (Textbook)

$$\underbrace{(f(x) - \bar{g}(x))^2}_{\text{bias}} + \underbrace{\text{Var}(x)}$$

BIAS VARIANCE DECOMPOSITION

Theorem: $E_D \left[(f(x) - g^D(x))^2 \right] = \underbrace{(f(x) - \bar{g}(x))^2}_{\text{expected error}} + \text{Var}(x)$

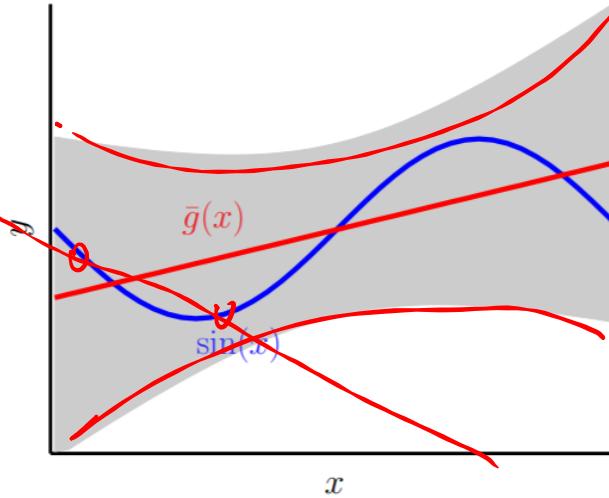
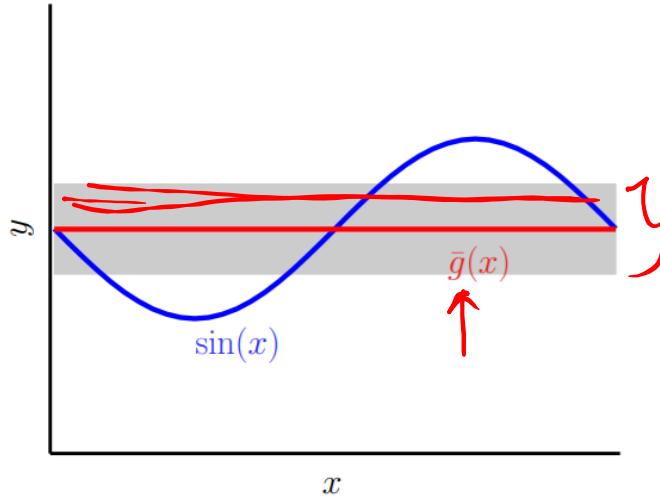
$$= \text{bias}(x) + \text{Var}(x)$$

$$E_X \left[E_D \left[(f(x) - g^D(x))^2 \right] \right] = E_X (\text{bias}(x)) + E_X (\text{Var}(x))$$

$h_0(x) = b \quad \} \text{Bias } \uparrow \text{ Variance } \downarrow$
 $h_{\theta}(x) = ax + b \quad \} \text{Bias } \downarrow \text{ Variance } \uparrow$
 $\uparrow \text{approx}$

generalization

Which one is better?



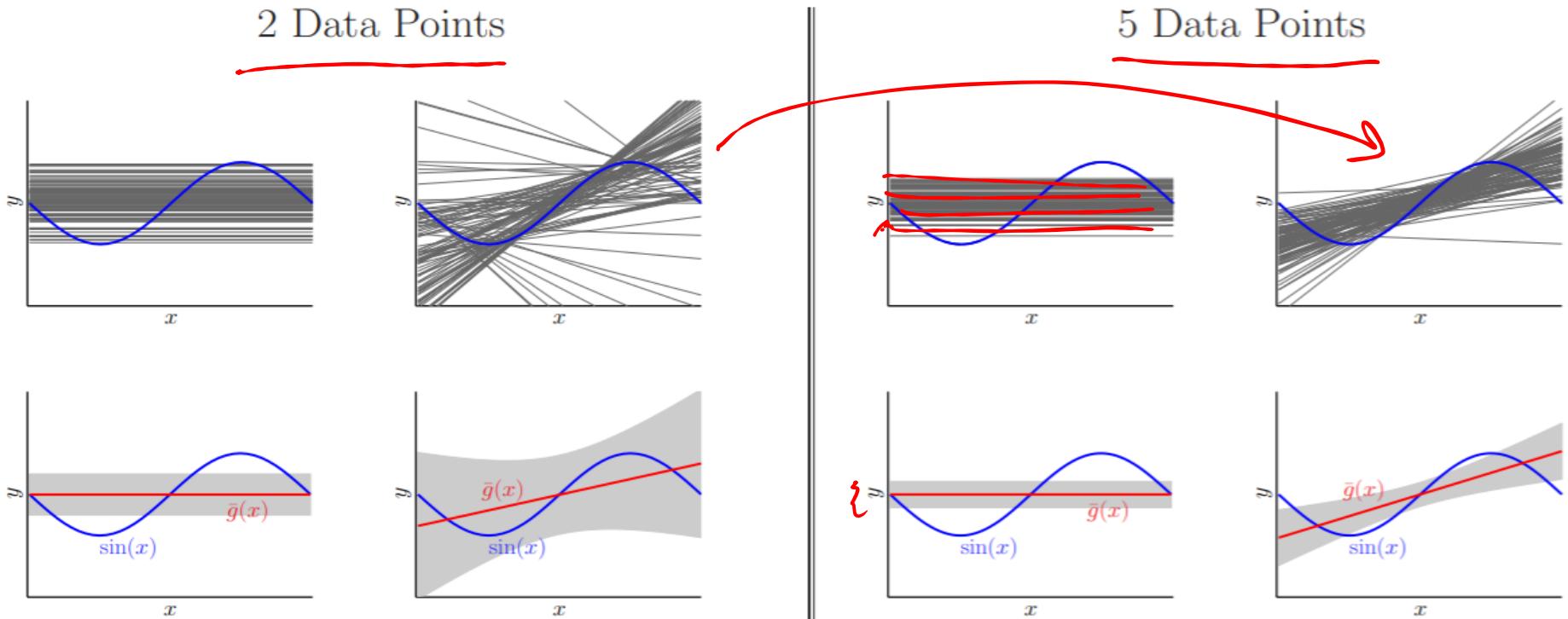
{

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

\checkmark

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

Data



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

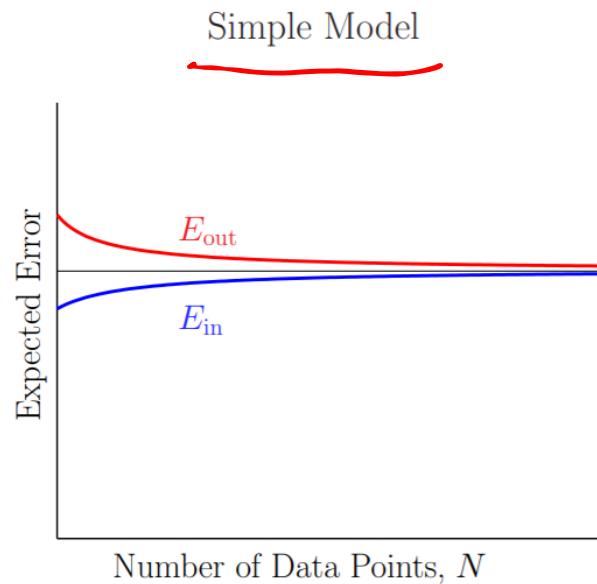
$$\begin{aligned} \mathcal{H}_1 & \\ \text{bias} &= 0.21; \\ \text{var} &= 1.69. \checkmark \\ \frac{E_{\text{out}}}{E_{\text{out}}} &= 1.90 \end{aligned}$$

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

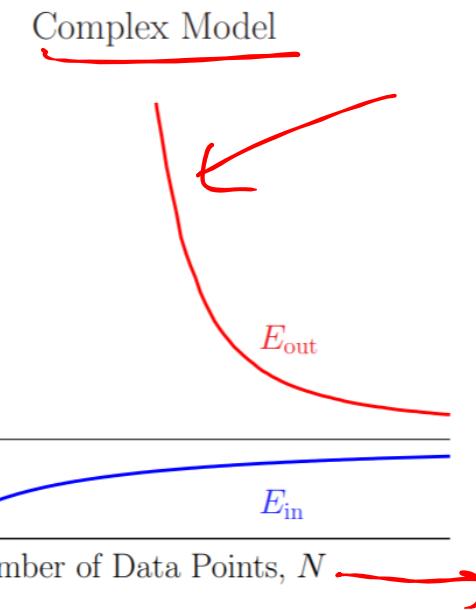
$$\begin{aligned} \mathcal{H}_1 & \\ \text{bias} &= 0.21; \\ \text{var} &= 0.21. \downarrow \\ \frac{E_{\text{out}}}{E_{\text{out}}} &= 0.42 \quad \checkmark \end{aligned}$$

Learning Curve

$$E_D [E_{out}(g^*)]$$

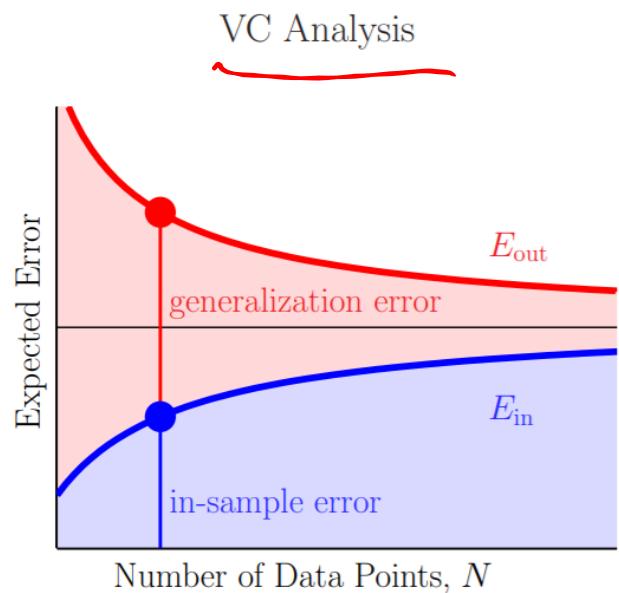


$$E_D [E_{in}(g^*)]$$

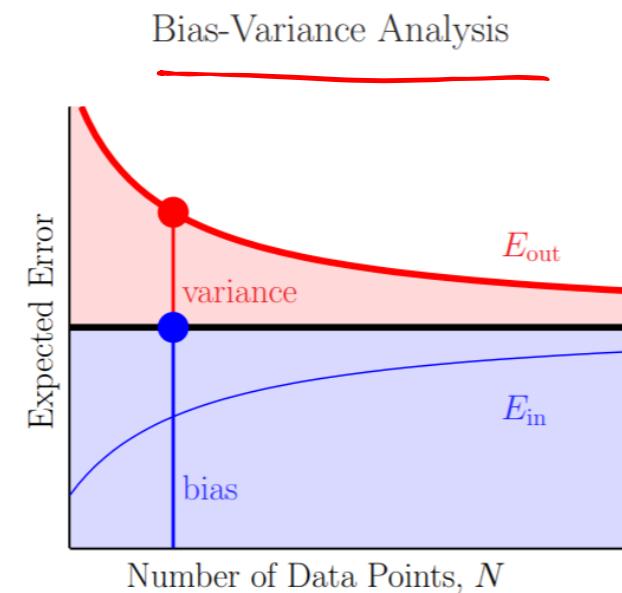


$$E_{out} = \mathbb{E}_{\mathbf{x}} [E_{out}(\mathbf{x})]$$

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

\bar{g}
 $\bar{g}(D)$

Thanks!