

CSE4088 Introduction to Machine Learning

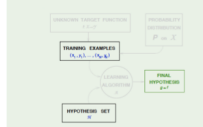
Bias-variance Tradeoff

Slides are adopted from lecture notes of Yaser Abu-Mostafa

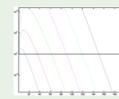
Review of last lecture

- VC dimension $d_{VC}(\mathcal{H})$
most points \mathcal{H} can shatter

- Scope of VC analysis



- Utility of VC dimension



Rule of thumb: $N \geq 10 d_{VC}$

- Generalization bound

$$E_{out} \leq E_{in} + \Omega$$

Ω is a function of N (number of samples), d_{VC} (VC dimension) and δ (prob. of error)

Outline

- Bias and variance
- Learning curves

Approximation-generalization tradeoff

Small E_{out} : good approximation of f out of sample.

More complex $\mathcal{H} \implies$ better chance of **approximating** f

Less complex $\mathcal{H} \implies$ better chance of **generalizing** out of sample

Ideal $\mathcal{H} = \{f\}$ winning lottery ticket ☺

Quantifying the tradeoff

VC analysis was one approach: $E_{out} \leq E_{in} + \Omega$

Bias-variance analysis is another: decomposing E_{out} into

1. How well \mathcal{H} can approximate f
2. How well we can zoom in on a good $h \in \mathcal{H}$

Applies to **real-valued targets** and uses **squared error**

Start with E_{out}

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

We want to release the dependency on the dataset by averaging out on D

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

(We can change the order of the expectations since the integrand is non-negative)

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

Now, let us focus on:

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

The average hypothesis

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

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

This will give us the correct value, in theory, if the noise is zero-mean.

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

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

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

Bias and variance

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

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

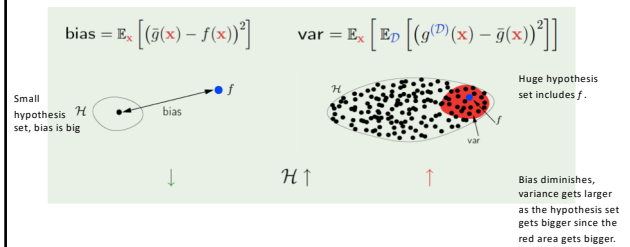
Var: How far is your hypothesis g from the average hypothesis (the best possible one)? It is the variance of the model due to the finite dataset.

Bias: How far is your best hypothesis from the target function? (You did your best by averaging out but it is still far from the target function)

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

Interpretation: You learn from one dataset and calculate your error E_{out} . You do this for other possible datasets and take the average. Lets say your error is 0.3. Then, 0.05 of it might come from bias and 0.25 of it might come from variance.

The tradeoff



Example: sine target

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

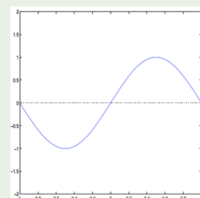
Only two training examples! $N = 2$

Two models used for learning:

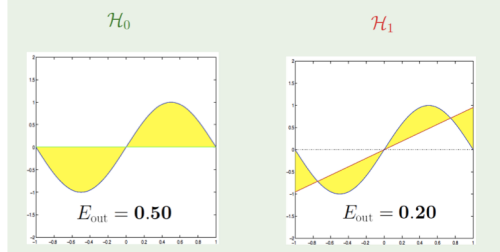
$$\mathcal{H}_0: h(x) = b$$

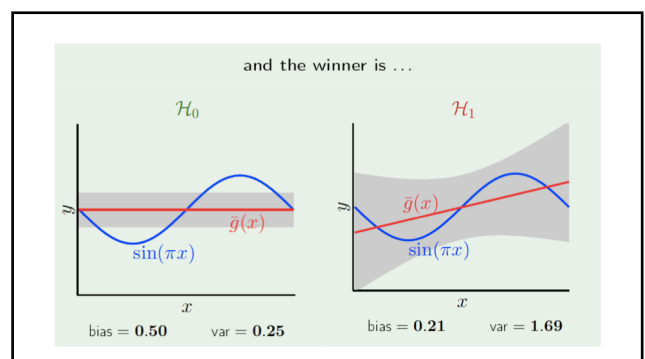
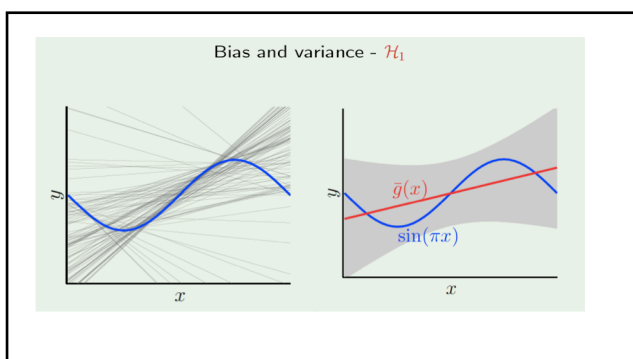
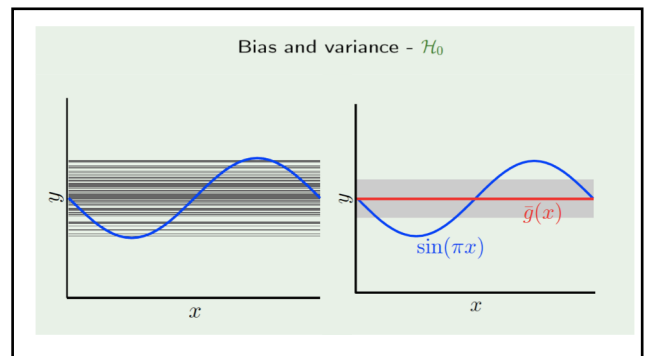
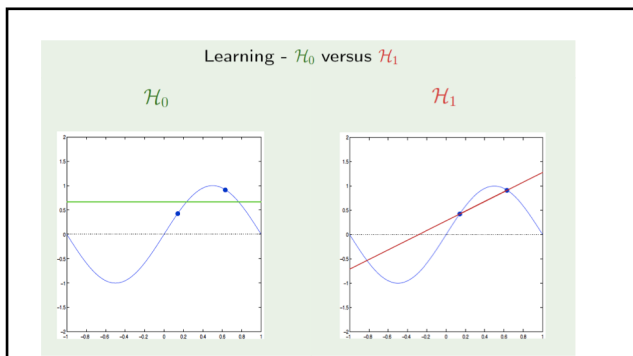
$$\mathcal{H}_1: h(x) = ax + b$$

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



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





Lesson learned

Match the 'model complexity'

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

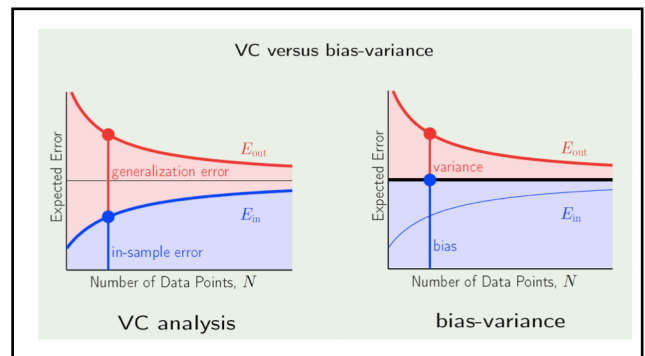
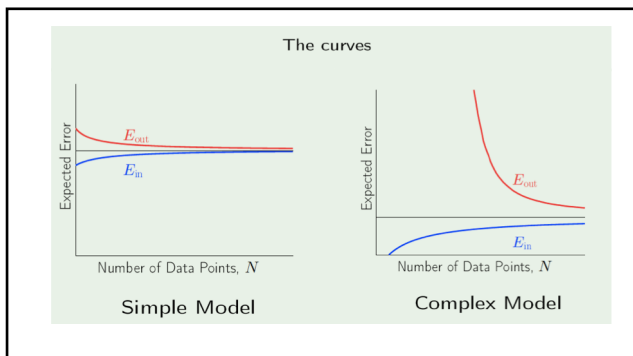
Expected E_{out} and E_{in}

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 ?



Linear regression case

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 = $X\mathbf{w} - \mathbf{y}'$

