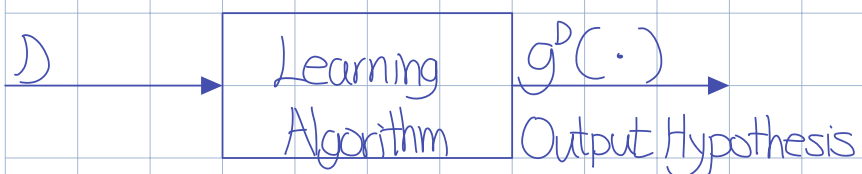


Bias Variance Tradeoff (Chapter 2)

Training set $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$



$D \sim P_D(\cdot)$ Distribution that generate D

let D_1, \dots, D_k be k independent dataset drawn from $P_D(\cdot)$



Average hypothesis

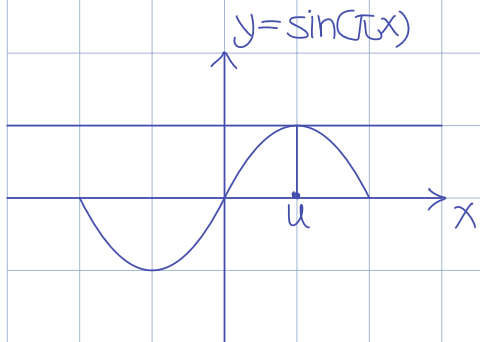
$$\frac{1}{k} \sum_{i=1}^k g^{D_i}(x)$$

$$\begin{aligned} \overline{g}(x) &= \lim_{k \rightarrow \infty} \frac{1}{k} \sum_{i=1}^k g^{D_i}(x) \\ &= \mathbb{E}_D[g^D(x)] \\ &= \int_D g^D(x) p(x) dD \end{aligned}$$

Example,

Unknown function $y = f(x) = \sin(\pi x)$

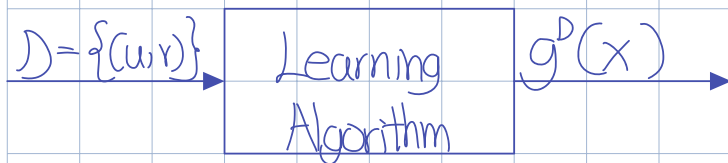
$d = 1$, $x \in \{-1, 1\}$



$D = \{(u, y)\}$ $N = 1$

$u \sim \text{unif}(-1, 1)$

$y = f(u) = \sin(\pi u)$

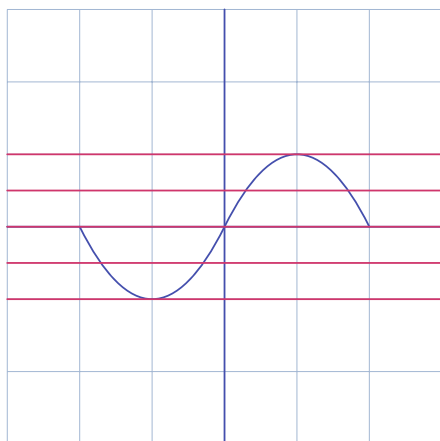


$\mathcal{H} \equiv$ Set of constant hypothesis

$g^D(x) = \text{Constant}$

Output Hypothesis

$g^u(x) = y = \sin(\pi u)$



$$\begin{aligned}\overline{g}(x) &= \mathbb{E}_D(g^D(x)) \\ &= \mathbb{E}_D(\sin(\pi u)) \\ &= \int_{-1}^1 \sin(\pi u) f(u) du = 0\end{aligned}$$

$\overline{g}(x)$ can be interpreted in 2 different ways

1. Output hypothesis in H with infinite amount of training data
2. Best approximation to $f(x)$ in H

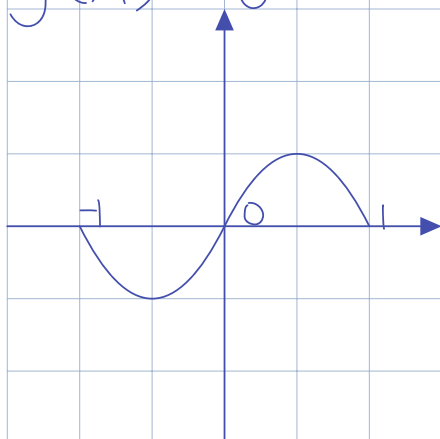
Bias Function

$$b(x) = (\overline{g}(x) - f(x))^2$$

$$f(x) = \sin(\pi x)$$

$$b(x) = (\sin(\pi x))^2$$

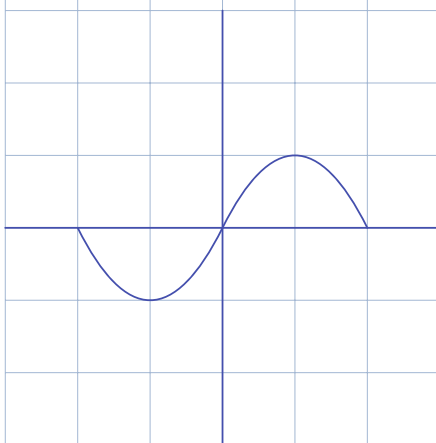
$$\overline{g}(x) = 0$$



Variance

$$\text{Var}(x) = \mathbb{E}_u(g^p(x) - \overline{g}(x))^2$$

$$g^p(x) \triangleq g^u(x) = \sin(\pi u)$$



$$\overline{g}(x) = 0$$

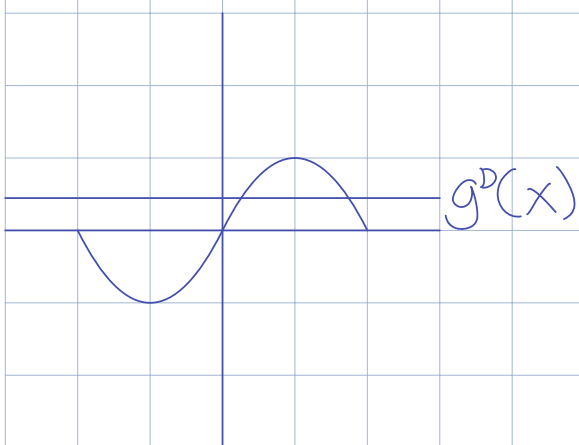
$$\text{Var}(x) = \mathbb{E}_u(\sin^2(\pi u))$$

$$= \int_{-1}^1 \sin^2(\pi u) p(u) du$$

$$= \frac{1}{2} \int_{-1}^1 \sin^2(\pi u) du$$

$$= \frac{1}{2}$$

pdf of u



Average test error

$$\Delta(x) = \mathbb{E}_u((g^p(x) - f(x))^2)$$

$$g^p(x) = g^u(x) = \sin(\pi u)$$

$$f(x) = \sin(\pi x)$$

$$\Delta(x) = \mathbb{E}_u[(\sin(\pi u) - \sin(\pi x))^2]$$

$$= E[\sin^2(\pi w) + \sin^2(\pi x) - 2\sin(\pi w)\sin(\pi x)]$$

$$= E[\sin^2(\pi w)] + E[\sin^2(\pi x)] - E[2\sin(\pi w)\sin(\pi x)]$$

$$= \frac{1}{2} + \sin^2(\pi x) - 2\sin(\pi x)E[\sin(\pi w)]$$

$$= \frac{1}{2} + \sin^2(\pi x)$$

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

Key Observation

$$\begin{aligned} \text{Average Test Error} &= \Delta(x) \\ &= \text{bias}(x) + \text{var}(x) \end{aligned}$$

Bias: How well can we approximate $y=f(x)$ using hypothesis class H , complexity of class

Variance: How much $g^D(x)$ fluctuates around $\bar{g}(x)$ as we vary D

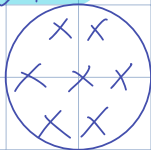
$f(x)$

Small H



Large bias small variance

Large H



$f(x)$

Large variance small bias

Recall

$$\Delta(x) = E_D((g^D(x) - f(x))^2)$$

$$b(x) = (\bar{g}(x) - f(x))^2$$

$$\text{Var}(x) = E_D((g^D(x) - \bar{g}(x))^2)$$

$$\Delta(x) = E_D[(g^D(x) - \bar{g}(x) + \bar{g}(x) - f(x))^2]$$

$$= E_D[(g^D(x) - \bar{g}(x))^2 + (\bar{g}(x) - f(x))^2 +$$

$$2(g^D(x) - \bar{g}(x))(\bar{g}(x) - f(x))]$$

$$= E_D[(g^D(x) - \bar{g}(x))^2] + E_D[(\bar{g}(x) - f(x))^2] +$$

$$2E_D[(g^D(x) - \bar{g}(x))(\bar{g}(x) - f(x))] \xrightarrow{\text{Zero}}$$

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

$$\begin{aligned}
 \overline{E_{out}} &= \int \Delta(x) p(x) dx \\
 &\quad \downarrow \\
 &\text{Average of Test Error} \\
 &= \int (\text{bias}(x) + \text{var}(x)) p(x) dx \\
 &= \int \text{bias}(x) p(x) dx + \int \text{var}(x) p(x) dx \\
 &= \overline{\text{bias}(x) + \text{var}(x)}
 \end{aligned}$$

Chapter 4, Validation



$$E_{val}(\bar{g}) = \frac{1}{K} \sum_{\substack{(x_n, y_n) \\ \in D_{val}}} e(\bar{g}(x_n), y_n)$$

How well does $E_{val}(\bar{g})$ approximate $E_{out}(\bar{g}) = \mathbb{E}_x(e(\bar{g}(x), y_n))$
Test Error

As $k \longrightarrow \infty$

$$E_{\text{in}}(\bar{g}) = E_{\text{out}}(\bar{g})$$

Note

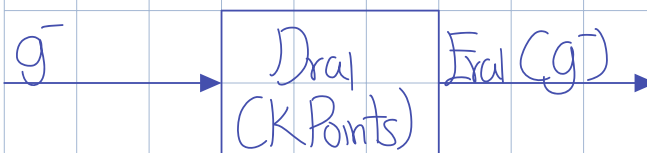
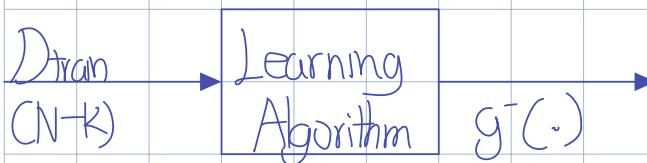
$E_{\text{in}}(\bar{g})$ is random

Since D_{in} is random $\sim P_D(\cdot)$

1. $E_{D_{\text{in}}} [E_{\text{in}}(\bar{g})] = E_{\text{out}}(\bar{g})$

2. With probability $\geq 1 - \delta$

$$E_{\text{out}}(\bar{g}) \leq E_{\text{in}}(\bar{g}) + \sqrt{\frac{1}{2k} \log \frac{2}{\delta}}$$



$$E_{\text{out}}(\bar{g}) = E_{D_{\text{in}}} [E_{\text{in}}(\bar{g})]$$

Not \bar{g} is fixed hypothesis wrt D_{in}

By Hoeffding Inequality,

$$\Pr(|E_{\text{in}}(\bar{g}) - E_{\text{out}}(\bar{g})| > \epsilon) \leq 2e^{-2k\epsilon^2}$$

$$\delta = 2e^{-2k\epsilon^2}$$

$$\epsilon = \sqrt{\frac{1}{2k} \log \frac{2}{\delta}}$$

With probability $\geq 1 - \delta$

$$|E_{\text{out}}(\bar{g}) - E_{\text{in}}(\bar{g})| < \sqrt{\frac{1}{2k} \log \frac{2}{\delta}}$$

$$E_{\text{out}}(\bar{g}) \leq E_{\text{in}}(\bar{g}) + \sqrt{\frac{1}{2k} \log \frac{2}{\delta}}$$

Proof of Part 1

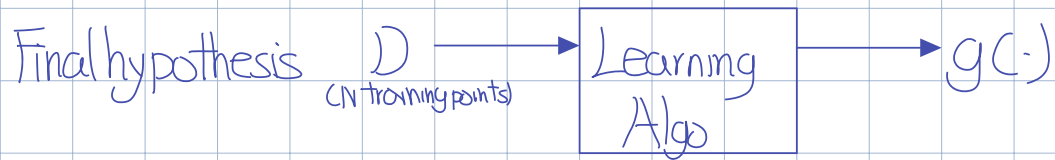
$$[E_{\text{in}}][E_{\text{in}}(\bar{g})] = E_{\text{out}}(\bar{g})$$

$$= E_{\mathbf{x}(1..k), \mathbf{X}^k} \left[\frac{1}{k} \sum_{\substack{(\mathbf{x}_n, y_n) \\ \in \text{D}_{\text{in}}}} e(\bar{g}(\mathbf{x}_n), y_n) \right]$$

$$= \frac{1}{k} \sum_{n=1}^k E_{\mathbf{x}_n} [e(\bar{g}(\mathbf{x}_n), y_n)]$$

$E_{\text{out}}(\bar{g})$

$\text{Var}(E_{\text{train}})$ decrease as $\frac{1}{k}$



$$E_{\text{out}}(g) \underset{\substack{N-k \\ \text{large}}}{\approx} E_{\text{out}}(g^-) \underset{\substack{k \\ \text{large}}}{\approx} E_{\text{train}}(g^-)$$

$$k = \frac{N}{5} \quad \text{Good}$$