

BIAS - VARIANCE TRADE-OFF I

1) Hypothesis space: \mathcal{H}

- \mathcal{H} is a compact subset of $C(X)$, equipped with the standard norm

$$\|f\|_{C(X)} := \|f\|_{\infty} := \sup_{\vec{x} \in X} |f(\vec{x})|$$

2) Empirical target function:

- h^* is not computable

- h^* may not reside in \mathcal{H}

\Rightarrow Thus the best we hope for is to find a target function \hat{h} in \mathcal{H} that is closest to h^*

$$\hat{h} := \operatorname{argmin}_{\mathcal{H}} \int_X (f(x) - h^*(x))^2 d\pi(\vec{x})$$

E+17.1. \hat{h} is also closest to y , that is

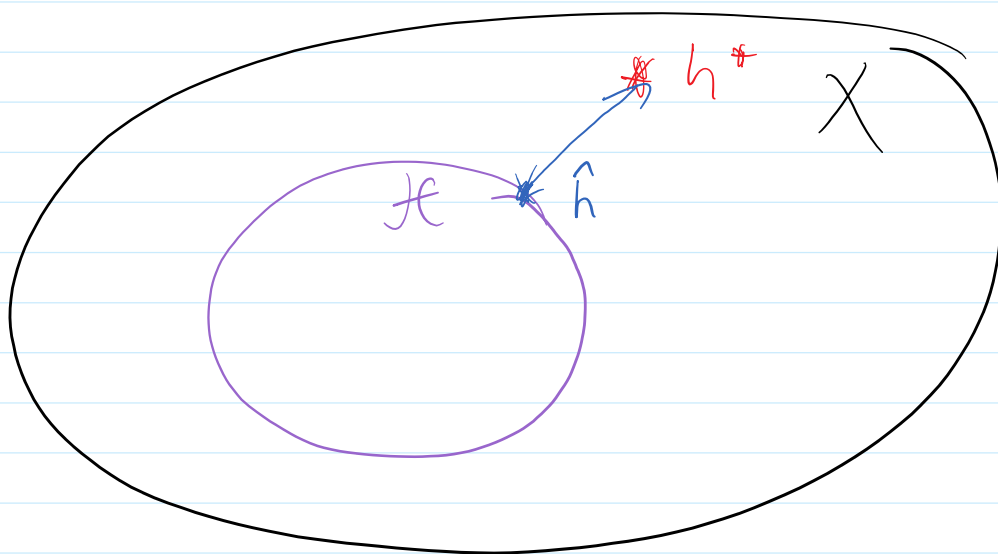
\hat{h} is also a minimizer of

$$\int_X (y - h)^2 d\pi(\vec{x})$$

$$\min_{\mathcal{H}} \int_{\mathcal{Z}} (f(\vec{x}) - y)^2 d\pi(\vec{x}, y)$$



$$\hat{h} := \operatorname{argmin}_{\mathcal{H}} \mathcal{R}(h)$$



- Again we know $\pi(\vec{x}, y)$ only through the training set $\Rightarrow \hat{h}$ is NOT computable. Thus we resort to empirical target function:

$$\hat{h}_N := \operatorname{argmin}_{f \in \mathcal{H}} \mathcal{R}_N(f) = \frac{1}{N} \sum_{i=1}^N (f(\vec{x}_i) - y_i)^2$$

EXISTENCE of \hat{h} , \hat{h}_N

Assumption: (M-Boundedness of the misfit on \mathcal{H})
for any $h \in \mathcal{H}$ and a.s. (a.e) in X
there holds:

$$|h(\vec{x}) - y| \leq M$$

Prop 17.1: Suppose \mathcal{H} is M-bounded. Then
 $\mathcal{R}, \mathcal{R}_N: \mathcal{H} \rightarrow \mathbb{R}$ are Lipschitz
continuous, i.e.

$$|\mathcal{R}(h_1) - \mathcal{R}(h_2)| \leq c \|h_1 - h_2\|_\infty$$

$$|\mathcal{R}_N(h_1) - \mathcal{R}_N(h_2)| \leq c \|h_1 - h_2\|_\infty$$

Proof:

$$|\mathcal{R}(h_1) - \mathcal{R}(h_2)| = \int_Z [(h_1 - y)^2 - (h_2 - y)^2] d\pi$$

$$\parallel$$
$$\left| \int_Z (h_1 + h_2 - 2y)(h_1 - h_2) d\pi \right|$$

// M-boundedness

$$2M \|h_1 - h_2\|_\infty$$

Cor. 17.1: Suppose f is M -bounded, then \hat{h} and \hat{h}_N exist.

Proof: direct consequence of Weierstrass theorem and Prop 17.1

BIAS - VARIANCE TRADE-OFF

- We are interested in the actual risk $R(\hat{h}_N)$. Let us start with

$$R(\hat{h}_N) = \underbrace{R(\hat{h}_N) - R(\hat{h})}_{S(\hat{h}_N)} + \underbrace{R(\hat{h})}_{B(\hat{h})}$$

Goal: is to bound $S(\hat{h}_N)$ and $B(\hat{h})$

- Let us first consider

$$\begin{aligned} S(\hat{h}_N) &= R_N(\hat{h}_N) - R_N(\hat{h}) + R(\hat{h}_N) - R_N(\hat{h}_N) + R_N(\hat{h}) - R(\hat{h}) \\ &\leq |R(\hat{h}_N) - R_N(\hat{h}_N)| + |R_N(\hat{h}) - R(\hat{h})| \end{aligned}$$

↑ because \hat{h}_N is a minimizer of R_N

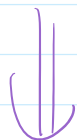
⇒ $J(\hat{h}_N)$ is bounded by sampling errors ⇒
hence the name. $J(\hat{h}_N)$ is also known as
the variances.

- Next let us look at

$$B(\hat{h}) := R(\hat{h}) \stackrel{??}{=} \int_X (\hat{h}(\vec{x}) - h^*(\vec{x}))^2 d\pi(\vec{x}) + O_{\pi}^2$$



The bias depends only the approximation
capability of the hypothesis space \mathcal{H} [closer
 h^* to \mathcal{H} the better the error].



the task of approximation theory
(NOT a focus of this class)

* Summary: we focus on estimating the
sampling / variance error in this class

* Bias - Variance trade-off

* Bias - Variance trade-off

- For fixed \mathcal{H} , the sampling error decreases as we increase the sample size N .
- For a fixed sample size N , enriching the hypothesis space \mathcal{H} reduces the bias

\Rightarrow the popular trade-off is to enlarge the hypothesis space as the sample size increases.