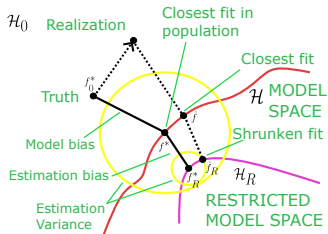


Introduction to Machine Learning

Regularization

Bias-variance Tradeoff



Learning goals

- Understand the bias-variance trade-off
- Know the definition of model bias, estimation bias, and estimation variance

BIAS-VARIANCE TRADEOFF

In this slide set, we will visualize the bias-variance trade-off.

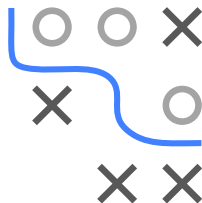
We consider a DGP \mathbb{P}_{xy} with $\mathcal{Y} \subset \mathbb{R}$ and the L2 loss L . We measure the distance between models $f : \mathcal{X} \rightarrow \mathbb{R}^g$ via

$$d(f, f') = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\mathbf{x}}} [L(f(\mathbf{x}), f'(\mathbf{x}))].$$

We define f_0^* as the risk minimizer such that

$$f_0^* \in \arg \min_{f \in \mathcal{H}_0} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{xy}} [L(y, f(\mathbf{x}))]$$

where $\mathcal{H}_0 = \{f : \mathcal{X} \rightarrow \mathbb{R} \mid d(\underline{0}, f) < \infty\}$ and $\underline{0} : \mathcal{X} \rightarrow \{0\}$.



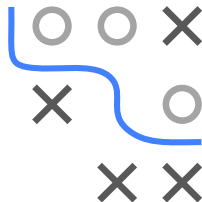
BIAS-VARIANCE TRADEOFF / 2

Our model space \mathcal{H} usually is a proper subset of \mathcal{H}_0 and in general $f_0^* \notin \mathcal{H}$.

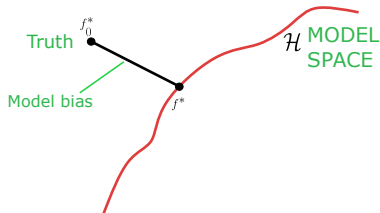
We define f^* as the risk minimizer in \mathcal{H} , i.e.,

$$f^* \in \arg \min_{f \in \mathcal{H}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{xy}} [L(f(\mathbf{x}, y))].$$

$f^* \in \mathcal{H}$ is closest to f_0^* , and we call $d(f_0^*, f^*)$ the model bias.

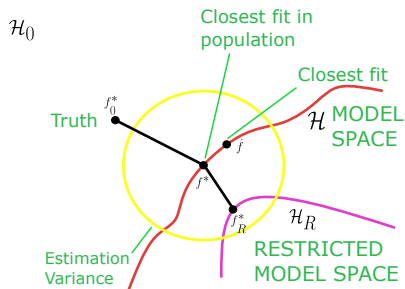


\mathcal{H}_0

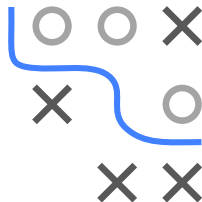


BIAS-VARIANCE TRADEOFF / 5

Let's assume that \hat{f} is an unbiased estimate of f^* (e.g., valid for linear regression), and we repeat the sampling process of \hat{f} .



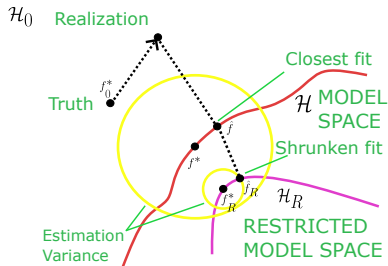
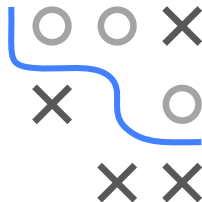
- We can measure the spread of sampled \hat{f} around f^* via $\delta = \text{Var}_{\mathcal{D}} [d(f^*, \hat{f})]$ which we call the estimation variance.
- We visualize this as a circle around f^* with radius δ .



BIAS-VARIANCE TRADEOFF / 6

We repeat the previous construction in the restricted model space \mathcal{H}_R and sample \hat{f}_R such that

$$\hat{f}_R \in \arg \min_{f \in \mathcal{H}_R} \sum_{i=1}^n L(y^{(i)}, \hat{f}(\mathbf{x}^{(i)})).$$



- We can measure the spread of sampled \hat{f}_R around f_R^* via $\delta = \text{Var}_{\mathcal{D}} [d(f_R^*, \hat{f}_R)]$ which we also call estimation variance.
- We observe that the increased bias results in a smaller estimation variance in \mathcal{H}_R compared to \mathcal{H} .