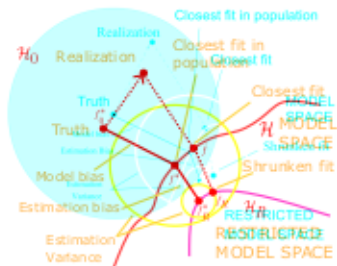


# Introduction to Machine Learning

## Regularization

## Bias-variance Tradeoff



### Learning goals

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



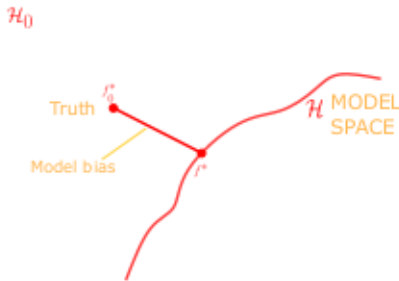
## BIAS-VARIANCE TRADEOFF / 2

Our model space  $\mathcal{H}$  is usually 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.,

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

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

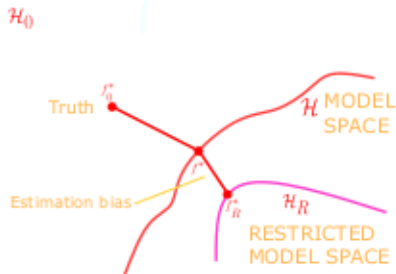


## BIAS-VARIANCE TRADEOFF / 3

By regularizing our model, we further restrict the model space so that  $\mathcal{H}_R$  is a proper subset of  $\mathcal{H}$ . We define  $f_R^*$  as the risk minimizer in  $\mathcal{H}_R$ , i.e.,

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

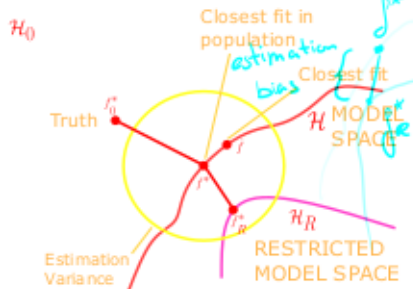
$f_R^* \in \mathcal{H}_R$  is closest to  $f_{\text{true}}$ , and we call  $d(f_R^*, f^*)$  the estimation bias.





## 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  $\hat{f}$ .



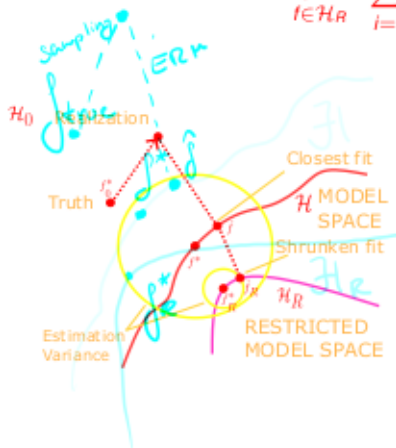
- We can measure the spread of sampled  $\hat{f}$  around  $f^*$  via  $\delta = \text{Var}_{\mathcal{D}} [d(f^*, \hat{f})]$  which we call the estimation variance.
- We visualize this as a circle around  $f^*$  with radius  $\delta$ .



# BIAS-VARIANCE TRADEOFF / 6

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

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



Note:

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



## BIAS-VARIANCE TRADEOFF / 7

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  $\delta$ .





## BIAS-VARIANCE TRADEOFF / 8

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

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



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