

# Statistical Learning

Wednesday, August 25, 2021 9:17 PM

## - Motivation

- We wish to determine the relationship between the inputs and outputs of a system.

↳ Assume such a relationship exists!

$Y \sim$  output (response, dependent variable)

$X \sim$  input (predictor, independent variable)

↳  $X \rightarrow (x_1, \dots, x_p)$  (features)  
dimension

$$Y = f(X) + \varepsilon$$

$\uparrow$   
systematic  
interaction  
of  $Y$  from  $X$

$\nwarrow$  error term  
independent of  $X$

- SL attempts to define and estimate  $f$ , we then attempt to evaluate the quality of our estimation!

\* In-Sample vs Out of Sample

## - Estimating $f$

- $f(X)$  is "ideal" or "best" representation we can have

↳  $\hat{f}(x)$

\* Why?

↳ Inference

↳ Prediction

- Prediction

$$\hat{y} = \hat{f}(x)$$

↳ "black box"

- Accuracy:  $\hat{y}$  as a prediction of  $y$

- Reducible Error  $\rightarrow$  Improve estimation of  $\hat{f}$

- Irreducible Error  $\rightarrow \epsilon$ , missing  $X$ 's

$$E[(y - \hat{y})^2] = E[(f(x) + \epsilon - \hat{f}(x))^2]$$

$$= \underbrace{(f(x) - \hat{f}(x))^2}_{\text{Reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible}}$$

↑

- Inference

• We estimate  $f$  in order to

- Understand association  $y \rightarrow x_1 \dots x_p$

- Which  $x_j$  are used, and considerations

\* Linearity  $\rightarrow$  Is  $y$  a linear function of  $x_1 \dots x_p$

- Parametric vs Non-Parametric Methods

\* Training Data  $\rightarrow$  Estimation of  $f$  using available sample data

- Parametric:

1. Assumption on function family / form of  $f$ :

$$\text{Linear: } f(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

2. Use training data to fit or draw  $\hat{f}(x)$

$$Y \approx \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

\* OLS Estimator

\* If  $f$  has a lot of  $\beta$ s  $\Rightarrow$  Flexible

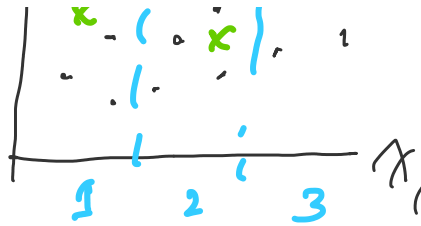
$\downarrow$   
overfitting  
(noise)

- Non-Parametric:

- No assumptions on structure of  $f$
- Estimate  $f$  as direct local/smoothed function of training data

$f(x)$  is piecewise defined over domain of  $X$

$x_4$  



- Accuracy

- Quality of fit/training : minimize error between our predicted values and actuals

- training (in-sample)
- test (out-of-sample)

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

↑  
mean
↑  
error
↑  
squared

\* sample  $x_i$  : iid

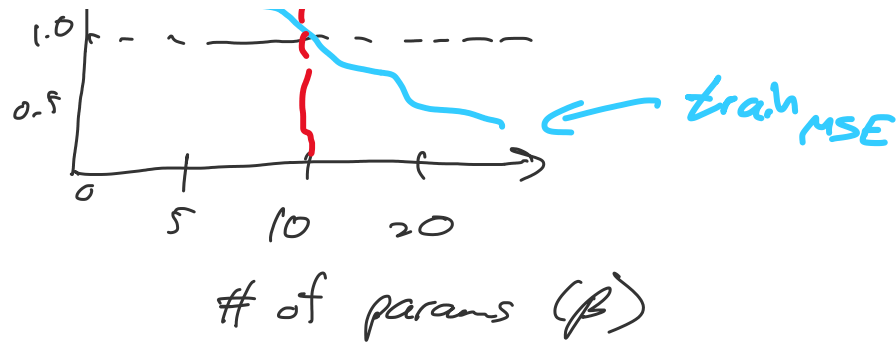
\* We have minimized error on training observations  $\{x_1, \dots, x_n\}$ , out-of-sample I get:

$(x_0, y_0)$  as test data

$$\text{Avg} (y_0 - \hat{f}(x_0))^2 \quad \leftarrow \text{in test!}$$

"generalize in test with  $\hat{f}(x)$ "





- Bias - Variance

• Expectations of Error

$$- E[\text{train MSE}] = \text{train MSE}$$

$$- E[\text{test MSE}] = ?$$

↳ Test Set separate from sample

↳ Test / Train split of sample

↳ k-Fold etc Cross-Validation (Holdouts)

\* for a given  $x_0$ ,  $E[\text{test MSE}]$

- variance of  $\hat{f}(x_0)$

- squared bias  $\hat{f}(x_0)$

- variance of error  $\varepsilon$

$$E[(y_0 - \hat{f}(x_0))^2] = \text{Var}(\hat{f}(x_0)) + (\text{Bias } \hat{f}(x_0))^2 + \text{Var}(\varepsilon)$$

↑  
Expected test MSE  
for any test  $x_0$

Variance  $\rightarrow$  change in  $\hat{f}$ 's params as  $x_0$  changes  
Bias  $\rightarrow$  Error in prediction is actual as  $x_0$  changes

# Trade-off: More flexible we make  $f$  in terms of # of parameters, the lower the bias, the higher the variance.