

The problem

Part I:

Goal: Given $\vec{x} = (x_1, \dots, x_p)$
 $y = \{ \text{response vector} \}$

Find a machine $M(\cdot)$ s.t.

$$\hat{y} = M(\vec{x})$$

Idea: Use training set $T = (\vec{x}_i, y_i)_{i=1}^N$
to train or build the machine
 $M(\cdot)$.

Part 2: Bias-Variance Problem

Assumption: $T = (X_i, Y_i)$

$$X_i \overset{\text{iid}}{\sim} f_X \quad Y_i | X_i \overset{\text{iid}}{\sim} f_{Y|X}$$

So we can think of T as a R.U.

Define the test set \vec{x}_0
and corresponding estimate $\hat{y}_0 = m(x_0)$

Part III: How good is my machine?

Expected - Prediction - Error = $EPE(\vec{x}_0)$

Assume we are in the setting

$$y = f(\vec{x}) + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

known \uparrow

So a given x_0 , $y_0 = f(x_0) + e_0$

known \nearrow \nwarrow stochastic

$$EPE(\vec{x}_0) \equiv E_{T, y_0} [(y_0 - \hat{y})^2]$$

$$= E_{y_0} [E_T [(y_0 - \hat{y}_i)^2]]$$

$$= \mathbb{E}_{Y_0} \left[\mathbb{E}_T \left[(Y_0 - \mathbb{E}_{Y_0}(Y_0)) + (\mathbb{E}_{Y_0}(Y_0) - \mathbb{E}_T(\hat{Y}_0)) + (\mathbb{E}_T(\hat{Y}_0) - \hat{Y}_0) \right]^2 \right]$$

$$= \mathbb{E}_{Y_0} \left[(Y_0 - \mathbb{E}_{Y_0}(Y_0))^2 \right] + (\mathbb{E}(Y_0) - \mathbb{E}_T(\hat{Y}_0))^2 + \mathbb{E}_T \left[(\mathbb{E}_T(\hat{Y}_0) - \hat{Y}_0)^2 \right]$$

+ expectations of cross terms
vanish

$$= \text{Var}(Y_0) + (\mathbb{E}(Y_0) - \mathbb{E}_T(\hat{Y}_0))^2 + \text{Var}_T(\hat{Y}_0)$$

$$= \sigma^2 + (f(x_0) - \mathbb{E}_T(\hat{Y}_0))^2 + \text{Var}_T(\hat{Y}_0)$$

↑
variance
of
target

↑
Bias²

↑
variance
of estimator

So

Thrm: $\mathbb{E}PE(\vec{x}_0) = \mathbb{E}_{T, Y_0}[(Y_0 - \hat{Y}_0)^2]$

$$= \text{Var}(Y) + \text{MSE}(\hat{Y}_0)$$

Linear Regression

Assume $f(\vec{x}) = \beta_0 + \sum_{j=1}^p \beta_j x_j = \vec{x}^T \vec{\beta}$

Goal: Find a method to come up with best $\hat{\beta}$.

$$X = \begin{bmatrix} \vec{x}_1^T \\ \vdots \\ \vec{x}_n^T \end{bmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix}$$

Strategy:

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\text{argmin}} \sum_{i=1}^n (Y_i - \beta^T x_i)^2$$

$$= \underset{\beta}{\operatorname{argmin}} (\gamma - X\beta)^T (\gamma - X\beta)$$

$$\frac{\partial}{\partial \beta} (\gamma - X\beta)^T (\gamma - X\beta) \equiv \nabla_{\beta} \operatorname{RSS}(\beta)$$

$$= \begin{pmatrix} \frac{d}{d\beta_0} \operatorname{RSS}(\beta) \\ \vdots \\ \frac{d}{d\beta_p} \operatorname{RSS}(\beta) \end{pmatrix} \stackrel{\text{set}}{=} \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$= -2X^T(\gamma - X\hat{\beta}) \stackrel{\text{set}}{=} 0$$

$$\Rightarrow X^T \gamma = X^T X \hat{\beta} \quad \text{normal equations}$$

$$\Rightarrow \text{If } X \text{ is full rank}$$

$$\hat{\beta}_{OLS} = (X^T X)^{-1} X^T Y$$

Properties

$$\hat{Y} = X\hat{\beta} = \underbrace{X(X^T X)^{-1} X^T}_{\textcircled{H}} Y = HY$$

orthogonal projection matrix.

onto $\text{col}(X)$

Preview of next lecture:

Ridge Regression

$$\hat{\beta}_R = \underset{\beta}{\text{argmin}} \underbrace{(Y - X\beta)^T (Y - X\beta)}_{\text{fit stat.}} + \underbrace{\lambda \beta^T \beta}_{\text{penalization}}$$

$$\hat{\beta}_R = (X^T X + \lambda I)^{-1} X^T y$$

"Shrinkage estimator"