

ESL- bias-variance decomposition notes

James Chuang

December 28, 2016

Bias-variance decomposition

From ESL pg. 23:

Suppose we have 1000 training examples x_i generated uniformly on $[-1, 1]^p$. Assume that the true relationship between X and Y is

$$Y = f(X) = e^{-8\|X\|^2}$$

, without any measurement error. We use the 1-nearest neighbor rule to predict y_0 at the test-point $x_0 = 0$. Denote the training set by τ . We can then compute the expected prediction error at x_0 for our procedure, averaging over all such samples of size 1000. Since the problem is deterministic, this is the mean squared error (MSE) for estimating $f(0)$:

$$\begin{aligned}\text{MSE}(x_0) &= \mathbb{E}_\tau [f(x_0) - \hat{y}_0]^2 \\ &= \mathbb{E}_\tau [\hat{y}_0 - f(x_0)]^2 \\ &= \mathbb{E}_\tau [\hat{y}_0 - \mathbb{E}_\tau [\hat{y}_0] + \mathbb{E}_\tau [\hat{y}_0] - \hat{y}_0]^2 \\ &= \mathbb{E}_\tau [y_0 - \mathbb{E}_\tau [\hat{y}_0]]^2 + 2\mathbb{E}_\tau [(\hat{y}_0 - \mathbb{E}_\tau [\hat{y}_0]) (\mathbb{E}_\tau [\hat{y}_0] - f(x_0))] + \mathbb{E}_\tau [\mathbb{E}_\tau [\hat{y}_0] - f(x_0)]^2 \\ &= \text{Var}_\tau(\hat{y}_0) + 2\mathbb{E}_\tau [(\hat{y}_0 - f(x_0)) (f(x_0) - f(x_0))] + \text{Bias}_\tau^2(\hat{y}_0) \\ &= \text{Var}_\tau(\hat{y}_0) + \text{Bias}_\tau^2(\hat{y}_0)\end{aligned}$$

Another way:

$$\begin{aligned}\text{MSE}(x_0) &= \mathbb{E}_\tau [f(x_0) - \hat{y}_0]^2 \\ &= \mathbb{E}_\tau [f^2(x_0) - 2f(x_0)\hat{y}_0 + \hat{y}_0^2] \\ &= \mathbb{E}_\tau [f^2(x_0)] - \mathbb{E}_\tau [2f(x_0)\hat{y}_0] + \mathbb{E}_\tau [\hat{y}_0^2] \\ &= \text{Var}_\tau(f(x_0)) + (\mathbb{E}_\tau [f(x_0)])^2 - 2f(x_0)\mathbb{E}_\tau [\hat{y}_0] + \text{Var}_\tau(\hat{y}_0) + (\mathbb{E}_\tau [\hat{y}_0])^2 \\ &= \text{Var}_\tau(f(x_0)) + \text{Var}_\tau(\hat{y}_0) + [\mathbb{E}_\tau(\hat{y}_0 - f(x_0))]^2 \\ &= \sigma^2 + \text{Var}_\tau(\hat{y}_0) + \text{Bias}^2(y_0)\end{aligned}$$