

1 Bias-Variance Trade-Off

When we are talking the concept of bias-variance trade off, try to image that one could repeatedly building a model multiple times based on different data (i.e. by resampling or gathering new data). Based on it, the following definitions are given.

Definition 1.1. Bias is the difference between the average prediction among our repeatedly built models and the actual value which we are trying to predict. Therefore,

$$\text{Bias}(\hat{f}(x_0)) = E[\hat{f}(x_0)] - f(x_0) \quad (1)$$

where x_0 is the sample that is trying to be predicted and $f(x_0)$ is the true value without considering the effect of noise ϵ (i.e. the true value corresponds to x_0 is y_0 and $y_0 = f(x_0) + \epsilon$).

Definition 1.2. Variance is defining how the model prediction varies for a given data point when repeatedly building models based on different sample data. Therefore,

$$\text{Var}(\hat{f}(x_0)) = E[\hat{f}(x_0) - E(\hat{f}(x_0))]^2 \quad (2)$$

where x_0 is the sample that is trying to be predicted and $f(x_0)$ is the true value without considering the effect of noise ϵ (i.e. the true value corresponds to x_0 is y_0 and $y_0 = f(x_0) + \epsilon$).

Based on definition 1.1 and 1.2, we know that model with high bias pays very little attention to the training data and oversimplifies the model, and model with high variance pays a lot of attention to training data and does not generalize on the data which it has not seen before.

Now, let's focusing on proving

$$E(y_0 - \hat{f}(x_0))^2 = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon) \quad (3)$$

Here, we also assume that $\epsilon \sim \mathcal{N}(0, \sigma^2)$, which means we assume that the noise is the gaussian white noise (i.e. mean = zero). For simplicity, the following derivative process uses y for representing y_0 , f for representing $f(x_0)$, and \hat{f} for representing $\hat{f}(x_0)$

$$\begin{aligned} E((y - \hat{f}))^2 &= E[(f + \epsilon - \hat{f})^2] \\ &= E\left[(f + \epsilon - \hat{f} + E(\hat{f}) - E(\hat{f}))^2\right] \\ &= E\left[((f - E(\hat{f})) - (\hat{f} - E(\hat{f})) + \epsilon)^2\right] \\ &= E\left[(f - E(\hat{f}))^2 + (\hat{f} - E(\hat{f}))^2 + \epsilon^2 - 2(f - E(\hat{f}))(\hat{f} - E(\hat{f}))\right. \\ &\quad \left.+ 2\epsilon(f - E(\hat{f})) - 2\epsilon(\hat{f} - E(\hat{f}))\right] \quad ^1 \\ &= E[(f - E(\hat{f}))^2] + E(\hat{f} - E(\hat{f}))^2 + E(\epsilon^2) \\ &\quad - 2E[(f - E(\hat{f}))(\hat{f} - E(\hat{f}))] + 2E[\epsilon(f - E(\hat{f}))] - 2E[\epsilon(\hat{f} - E(\hat{f}))] \end{aligned} \quad (4)$$

¹ $\forall a, b, c \in R, (a + b + c)^2 = a^2 + b^2 + c^2 + 2ab + 2ac + 2bc.$

Since f is deterministic, $f - E(\hat{f})$ is a constant. Thus,

$$\begin{aligned}
 E((y - \hat{f}))^2 &= E[(f - E(\hat{f}))^2] + E(\hat{f} - E(\hat{f}))^2 + E(\epsilon^2) \\
 &\quad - 2E[(f - E(\hat{f}))(\hat{f} - E(\hat{f}))] + 2E[\epsilon(f - E(\hat{f}))] - 2E[\epsilon(\hat{f} - E(\hat{f}))] \\
 &= (f - E(\hat{f}))^2 + E(\hat{f} - E(\hat{f}))^2 + E(\epsilon^2) \\
 &\quad - 2(f - E(\hat{f}))E(\hat{f} - E(\hat{f})) + 2(f - E(\hat{f}))E(\epsilon) - 2E(\epsilon)E(\hat{f} - E(\hat{f})) \quad ^2
 \end{aligned} \tag{5}$$

Based on definition 1.1 and 1.2, we know that the term $(f - E(\hat{f}))^2$ in (5) is $[\text{Bias}(\hat{f})]^2$, and the term $E(\hat{f} - E(\hat{f}))^2$ is $\text{Var}(\hat{f})$. Also, based on our assumption that $E(\epsilon) = 0$, we know that $\text{Var}(\epsilon) = E(\epsilon^2)$. Therefore,

$$E((y - \hat{f}))^2 = \text{Var}(\hat{f}) + [\text{Bias}(\hat{f})]^2 + \text{Var}(\epsilon) - 2(f - E(\hat{f}))E(\hat{f} - E(\hat{f})) \tag{6}$$

Since $E(\hat{f} - E(\hat{f})) = E(\hat{f}) - E(\hat{f}) = 0$, the final conclusion would be

$$E((y - \hat{f}))^2 = \text{Var}(\hat{f}) + [\text{Bias}(\hat{f})]^2 + \text{Var}(\epsilon) \tag{7}$$

2 References

[1] S, Fortmann-Roe, Understanding the Bias-Variance Tradeoff, June. 2012. Accessed on: <http://scott.fortmann-roe.com/docs/BiasVariance.html>

² Some background knowledge that should be known: 1. $\forall a, b \in R, E(a + bX) = a + bE(X)$; 2. If X, Y are two independent variables, then $E(XY) = EX \cdot EY$. That is why $E[\epsilon(\hat{f} - E(\hat{f}))] = E(\epsilon)E(\hat{f} - E(\hat{f}))$.