

Consider the regression setting, where we would like to design a learning algorithm that given a dataset \mathcal{D} consisting of N i.i.d. samples: $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ ($\mathbf{x}_i \in \mathbb{R}^d, y \in \mathbb{R}$), produces a predictor function (or hypothesis) $h : \mathbb{R}^d \rightarrow \mathbb{R}$. In this section, we analyze the generalization error of a learning algorithm via the so-called *Bias-Variance Decomposition* (or *Tradeoff*).

To make things concrete for our analysis, assume we sample from the joint distribution $\mathbf{x}, y \sim P(\mathbf{X}, Y)$ *, and therefore $\mathcal{D} \sim P^N(\mathbf{X}, Y)$ (i.e., independently sample N times). Given a dataset \mathcal{D} , our learning algorithm produces a function $h_{\mathcal{D}} : \mathbb{R}^d \rightarrow \mathbb{R}$ (note that this is dependent on \mathcal{D}). And the quantity we are interested in, the expected test error is:

$$\mathbb{E}_{\mathcal{D} \sim P^N, (\mathbf{x}, y) \sim P} [(h_{\mathcal{D}}(\mathbf{x}) - y)^2]. \quad (1)$$

Before we begin the analysis, let us first define a few useful notations:

- (*Expected Label*) Denote the conditional expectation of y given \mathbf{x} as

$$\bar{y}(\mathbf{x}) = \mathbb{E}_{Y|X=\mathbf{x}}[Y|\mathbf{X} = \mathbf{x}] = \int_y y P(y|\mathbf{X} = \mathbf{x}) dy.$$

- (*Expected Prediction*) Fixing our learning algorithm, denote the expected prediction given \mathbf{x} over the dataset \mathcal{D} distribution as

$$\bar{h}(\mathbf{x}) = \mathbb{E}_{\mathcal{D} \sim P^N}[h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) P(\mathcal{D}) d\mathcal{D}. \quad (2)$$

We are now ready to derive the *Bias-Variance Decomposition* of (1):

$$\begin{aligned} & \mathbb{E}_{\mathcal{D} \sim P^N, (\mathbf{x}, y) \sim P} [(h_{\mathcal{D}}(\mathbf{x}) - y)^2] \\ &= \mathbb{E}_{\mathcal{D} \sim P^N, (\mathbf{x}, y) \sim P} [((h_{\mathcal{D}}(\mathbf{x}) - \bar{h}(\mathbf{x})) + (\bar{h}(\mathbf{x}) - y))^2] \\ &= \mathbb{E}_{\mathcal{D} \sim P^N, (\mathbf{x}, y) \sim P} [(h_{\mathcal{D}}(\mathbf{x}) - \bar{h}(\mathbf{x}))^2] \\ & \quad + 2\mathbb{E}_{\mathcal{D} \sim P^N, (\mathbf{x}, y) \sim P} [(h_{\mathcal{D}}(\mathbf{x}) - \bar{h}(\mathbf{x}))(\bar{h}(\mathbf{x}) - y)] \\ & \quad + \mathbb{E}_{(\mathbf{x}, y) \sim P} [(\bar{h}(\mathbf{x}) - y)^2]. \end{aligned}$$

Observe the middle term is 0 since

$$\begin{aligned} & \mathbb{E}_{\mathcal{D} \sim P^N, (\mathbf{x}, y) \sim P} [(h_{\mathcal{D}}(\mathbf{x}) - \bar{h}(\mathbf{x}))(\bar{h}(\mathbf{x}) - y)] \\ &= \mathbb{E}_{(\mathbf{x}, y) \sim P} [\mathbb{E}_{\mathcal{D} \sim P^N} [(h_{\mathcal{D}}(\mathbf{x}) - \bar{h}(\mathbf{x}))(\bar{h}(\mathbf{x}) - y) | (\mathbf{X}, Y)]] \\ & \quad \text{(Law of Total Expectation[†] and independence of } \mathcal{D} \text{ and } (\mathbf{x}, y)) \\ &= \mathbb{E}_{(\mathbf{x}, y) \sim P} [(\bar{h}(\mathbf{x}) - y) \cdot \mathbb{E}_{\mathcal{D} \sim P^N} [(h_{\mathcal{D}}(\mathbf{x}) - \bar{h}(\mathbf{x})) | (\mathbf{X}, Y)]] \\ &= \mathbb{E}_{(\mathbf{x}, y) \sim P} [(\bar{h}(\mathbf{x}) - y) \cdot 0] \quad \text{(since (2))} \\ &= 0. \end{aligned}$$

Hence,

$$\mathbb{E}_{\mathcal{D} \sim P^N, (\mathbf{x}, y) \sim P} [(h_{\mathcal{D}}(\mathbf{x}) - y)^2] = \underbrace{\mathbb{E}_{\mathcal{D} \sim P^N, (\mathbf{x}, y) \sim P} [(h_{\mathcal{D}}(\mathbf{x}) - \bar{h}(\mathbf{x}))^2]}_{\text{Variance}} + \mathbb{E}_{(\mathbf{x}, y) \sim P} [(\bar{h}(\mathbf{x}) - y)^2]. \quad (3)$$

*In an slight abuse of notation, we here used P to denote the density function.

[†]This is a special application of the *Law of Total Expectation*, which many of you probably have already used in HW0.

We can easily derive this from definition: $\mathbb{E}_{X,Y}[g(X,Y)] = \int_x \int_y g(x,y)p(x,y)dydx = \int_x \int_y g(x,y)p(y|x)p(x)dydx = \int_x \left(\int_y g(x,y)p(y|x)dy \right) p(x)dx = \int_x \mathbb{E}_{Y|X=x}[g(X,Y)|X=x]p(x)dx = \mathbb{E}_X[\mathbb{E}_{Y|X}[g(X,Y)|X]].$

Note that at this point, we already have the Variance term: $\mathbb{E}_{\mathcal{D} \sim P^N, (\mathbf{x}, y) \sim P} \left[(h_{\mathcal{D}}(\mathbf{x}) - \bar{h}(\mathbf{x}))^2 \right]$, which expresses, over all possible dataset \mathcal{D} , how much does the prediction generated by our algorithm deviate from the expected prediction in expectation. The higher this quantity is, the bigger the change in prediction when we vary our dataset. In particular, when *overfitting* happens, our learning algorithm produces prediction very finetuned for every dataset, thus resulting in a high deviation from the expected prediction.

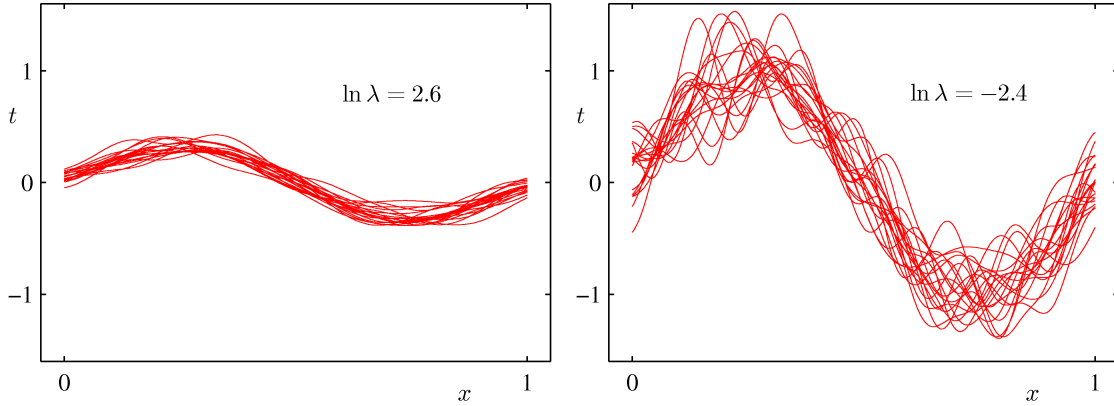


Figure 1: Plot of 20 of 100 predictor functions produced by the linear regression with regularization $\lambda \in \{\exp(2.6), \exp(-2.4)\}$ learning algorithm, trained with sample $N = 25$ over 100 datasets. Examples of low/high variance learning algorithms. Taken from [1].

We can apply the same trick again for the second term in (3) and complete the *Bias-Variance Decomposition*:

$$\begin{aligned}
 & \mathbb{E}_{(\mathbf{x}, y) \sim P} \left[(\bar{h}(\mathbf{x}) - y)^2 \right] \\
 &= \mathbb{E}_{(\mathbf{x}, y) \sim P} \left[((\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) + (\bar{y}(\mathbf{x}) - y))^2 \right] \\
 &= \mathbb{E}_{(\mathbf{x}, y) \sim P} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))^2 \right] + \mathbb{E}_{(\mathbf{x}, y) \sim P} \left[(\bar{y}(\mathbf{x}) - y)^2 \right] \\
 &\quad + 2\mathbb{E}_{(\mathbf{x}, y) \sim P} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))(\bar{y}(\mathbf{x}) - y) \right] \\
 &= \underbrace{\mathbb{E}_{(\mathbf{x}, y) \sim P} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))^2 \right]}_{\text{Bias}^2} + \underbrace{\mathbb{E}_{(\mathbf{x}, y) \sim P} \left[(\bar{y}(\mathbf{x}) - y)^2 \right]}_{\text{Noise}}.
 \end{aligned} \tag{4}$$

We can show the last term in (4) is 0 in a similar way as we did before:

$$\begin{aligned}
 & \mathbb{E}_{(\mathbf{x}, y) \sim P} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))(\bar{y}(\mathbf{x}) - y) \right] \\
 &= \mathbb{E}_{\mathbf{x} \sim \mathbf{X}} \left[\mathbb{E}_{y \sim Y | \mathbf{X} = \mathbf{x}} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))(\bar{y}(\mathbf{x}) - y) \mid \mathbf{X} = \mathbf{x} \right] \right] && \text{(Law of Total Expectation)} \\
 &= \mathbb{E}_{\mathbf{x} \sim \mathbf{X}} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) \mathbb{E}_{y \sim Y | \mathbf{X} = \mathbf{x}} \left[(\bar{y}(\mathbf{x}) - y) \mid \mathbf{X} = \mathbf{x} \right] \right] \\
 &= \mathbb{E}_{\mathbf{x} \sim \mathbf{X}} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) (\bar{y}(\mathbf{x}) - \mathbb{E}_{y \sim Y | \mathbf{X} = \mathbf{x}} [y \mid \mathbf{X} = \mathbf{x}]) \right] \\
 &= \mathbb{E}_{\mathbf{x} \sim \mathbf{X}} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) (\bar{y}(\mathbf{x}) - \bar{y}(\mathbf{x})) \right] && \text{(Definition of } \bar{y}(\mathbf{x}) \text{)} \\
 &= \mathbb{E}_{\mathbf{x} \sim \mathbf{X}} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) \cdot 0 \right] \\
 &= 0.
 \end{aligned}$$

The squared bias term $\mathbb{E}_{(x,y) \sim P} \left[(\bar{h}(x) - \bar{y}(x))^2 \right]$ is the expected squared difference between the expected prediction $\bar{h}(x)$ (over all datasets) and the expected label $\bar{y}(x)$ (over all datasets). In other words, this term represents the extent to which the average prediction over all data sets differs from the ground truth function. The higher the Bias², the further our learning algorithm like to "bias" away from the ground truth function (therefore inherently unable to approximate the function).

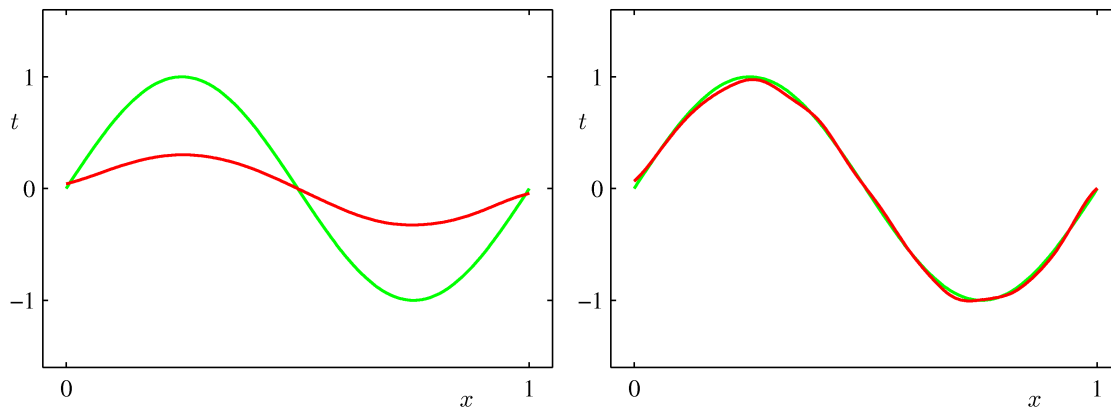


Figure 2: Plot of *average* of 100 predictor functions produced with the same setting as in Figure 1. Examples of high (left)/low (right) bias learning algorithms. Taken from [1].

Finally, the Noise term $\mathbb{E}_{(x,y) \sim P} [(\bar{y}(x) - y)^2]$ is the expected (over x) variance of sampled label y . It is intrinsic to the data and sampling process and cannot be erased by any learning algorithm (notice that this term is not involved with \mathcal{D} in any way). We have now completed our derivation of the *Bias-Variance Decomposition*, summarized by the following equation:

$$\mathbb{E}_{\mathcal{D} \sim P^N, (x,y) \sim P} [(h_{\mathcal{D}}(\mathbf{x}) - y)^2] = \underbrace{\mathbb{E}_{\mathcal{D} \sim P^N, (x,y) \sim P} [(h_{\mathcal{D}}(\mathbf{x}) - \bar{h}(\mathbf{x}))^2]}_{\text{Variance}} + \underbrace{\mathbb{E}_{(x,y) \sim P} [(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))^2]}_{\text{Bias}^2} + \underbrace{\mathbb{E}_{(x,y) \sim P} [(\bar{y}(\mathbf{x}) - y)^2]}_{\text{Noise}}.$$

References

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag, 2006. ISBN: 0387310738.