

Another view on generalization

We have formalized the problem of supervised learning as finding a function h in a given set \mathcal{H} that minimizes the risk $R(h)$. In the context of classification we hope to approximate the Bayes classifier, while in the context of regression we hope to approximate some true underlying function. We have already seen that the choice of \mathcal{H} must strike a delicate tradeoff between two desirable characteristics:

- a more complex \mathcal{H} leads to better chance of *approximating* ideal classifier/function;
- a less complex \mathcal{H} leads to better chance of *generalizing* to unseen data.

The problem with a more complex \mathcal{H} is that, while it can do a good job of approximating the ideal classifier/function, in the absence of a sufficient amount of data it is hard to determine the *correct* function from data and we run the risk of *overfitting*.

In the context of classification, we have already seen that the tradeoff can be precisely quantified in terms of the *VC generalization bound*, which takes the form

$$R(h) \leq \widehat{R}_n(h) + \epsilon(\mathcal{H}, n) \text{ with high probability.}$$

Here we develop an alternative method to quantify the tradeoff called the *bias-variance decomposition* which takes the form

$$R(h) \approx \text{noise} + \text{bias}^2 + \text{variance.}$$

Here, the *noise* represents fundamental/unavoidable error inherent in the problem, the *bias* captures how well \mathcal{H} can approximate the optimal h^* , while the *variance* captures how likely we are to pick a good $h \in \mathcal{H}$. This approach generalizes more easily to regression than the VC dimension approach developed for classification.

The bias-variance decomposition

We formalize the bias-variance tradeoff assuming the following:

- $h^* : \mathbb{R}^d \rightarrow \mathbb{R}$ is the unknown target function that we are trying to learn;
- $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ is the dataset, where (\mathbf{x}_i, y_i) are independent and identically distributed; specifically, $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i = f(\mathbf{x}_i) + \varepsilon_i \in \mathbb{R}$, where ε_i is a zero-mean noise random variable independent of \mathbf{x}_i with variance σ_ε^2 (for instance $\varepsilon_i \sim \mathcal{N}(0, \sigma_\varepsilon^2)$);
- $h_{\mathcal{D}} : \mathbb{R}^d \rightarrow \mathbb{R}$ is our choice of function in \mathcal{H} , selected using \mathcal{D} ;
- The performance of $h_{\mathcal{D}}$ is measured in terms of the mean squared error $R(h_{\mathcal{D}}) = \mathbb{E}_{XY} (h_{\mathcal{D}}(X) - Y)^2$;

Note that the random variables (X, Y) denote the data at *testing* and should not be confused with the random variables in \mathcal{D} representing the *training* data, and that since \mathcal{D} is itself random, the metric $R(h_{\mathcal{D}})$ is random and we will ultimately be interested in $\mathbb{E}_{\mathcal{D}} [R(h_{\mathcal{D}})]$.

Lemma 1 (Bias-variance decomposition)

$$\mathbb{E}_{\mathcal{D}} [R(h_{\mathcal{D}})] = \sigma_\varepsilon^2 + \mathbb{E}_X [\text{var}(h_{\mathcal{D}}(X)) | X] + \mathbb{E}_X [\text{bias}(h_{\mathcal{D}}(X))^2 | X]$$

with

$$\text{var}(h_{\mathcal{D}}(X)) = \mathbb{E}_{\mathcal{D}} \left[(h_{\mathcal{D}}(X) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(X))^2 \right]$$

$$\text{bias}(h_{\mathcal{D}}(X)) = \mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(X)] - h^*(X)$$

Proof To simplify notation, we set $\bar{h}(X) = \mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(X)]$. Then,

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}} [R(h_{\mathcal{D}})] &= \mathbb{E}_{\mathcal{D}} [\mathbb{E}_{XY} [(h_{\mathcal{D}}(X) - Y)^2]] \\
&= \mathbb{E}_{\mathcal{D}} [\mathbb{E}_{X\varepsilon} [(h_{\mathcal{D}}(X) - h^*(X) - \varepsilon)^2]] \\
&= \mathbb{E}_{\mathcal{D}} [\mathbb{E}_{X\varepsilon} [(h_{\mathcal{D}}(X) - \bar{h}(X) + \bar{h}(X) - h^*(X) - \varepsilon)^2]] \\
&= \mathbb{E}_{\mathcal{D}} \mathbb{E}_X \mathbb{E}_{\varepsilon} [(h_{\mathcal{D}}(X) - \bar{h}(X))^2 + (\bar{h}(X) - h^*(X))^2 + \varepsilon^2 \\
&\quad + 2(h_{\mathcal{D}}(X) - \bar{h}(X))(\bar{h}(X) - h^*(X)) \\
&\quad - 2(h_{\mathcal{D}}(X) - \bar{h}(X))\varepsilon - 2(\bar{h}(X) - h^*(X))\varepsilon]
\end{aligned}$$

Note that in the final equality above we have used the fact that \mathcal{D} , X , and ε are independent. Notice that

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}} \mathbb{E}_X \mathbb{E}_{\varepsilon} [(h_{\mathcal{D}}(X) - \bar{h}(X))^2] &= \mathbb{E}_X [\text{var}(h_{\mathcal{D}}(X)|X)] \\
\mathbb{E}_{\mathcal{D}} \mathbb{E}_X \mathbb{E}_{\varepsilon} [(\bar{h}(X) - h^*(X))^2] &= \mathbb{E}_X [\text{bias}(h_{\mathcal{D}}(X))^2] \\
\mathbb{E}_{\mathcal{D}} \mathbb{E}_X \mathbb{E}_{\varepsilon} [\varepsilon^2] &= \sigma_{\varepsilon}^2.
\end{aligned}$$

The last three terms turn out to be zero since

$$\begin{aligned}
&\mathbb{E}_{\mathcal{D}} \mathbb{E}_X [(h_{\mathcal{D}}(X) - \bar{h}(X))(\bar{h}(X) - h^*(X))] \\
&= \mathbb{E}_X \mathbb{E}_{\mathcal{D}} [(h_{\mathcal{D}}(X) - \bar{h}(X))(\bar{h}(X) - h^*(X))] \\
&= \mathbb{E}_X [(\mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(X)] - \bar{h}(X))(\bar{h}(X) - h^*(X))] \\
&= 0
\end{aligned}$$

and

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}} \mathbb{E}_X \mathbb{E}_{\varepsilon} [(h_{\mathcal{D}}(X) - \bar{h}(X))\varepsilon] &= \mathbb{E}_X [\mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(X) - \bar{h}(X)]] \mathbb{E}_{\varepsilon}[\varepsilon] = 0 \\
\mathbb{E}_{\mathcal{D}} \mathbb{E}_X \mathbb{E}_{\varepsilon} [(\bar{h}(X) - h^*(X))\varepsilon] &= \mathbb{E}_X [\bar{h}(X) - h^*(X)] \mathbb{E}_{\varepsilon}[\varepsilon] = 0.
\end{aligned}$$

■

Example

We now consider a concrete example. Suppose that $h^*(X) = \sin(\pi X)$ where $X \sim \text{Uniform}[-1, 1]$ and that we are given $n = 2$ noise-free training examples. We will consider two possible hypothesis sets:

- $\mathcal{H}_0 = \{h : h(x) = b, b \in \mathbb{R}\},$
- $\mathcal{H}_1 = \{h : h(x) = ax + b, a, b \in \mathbb{R}\}.$

Since the observations are noise-free (meaning each observation is of the form $Y = \sin(\pi X)$ with no ε term), the noise component of $R(h_{\mathcal{D}})$ is zero and all we need to compute are the bias and variance terms for each hypothesis set. We will perform these computations for each hypothesis set in turn.

Bias-variance decomposition for \mathcal{H}_0

We will begin, for reference, by computing

$$h^\# = \arg \min_{h \in \mathcal{H}_0} \mathbb{E}_X [(h - h^*(X))^2].$$

Since any $h \in \mathcal{H}_0$ is just a constant, this reduces to determining

$$b^* = \arg \min_{b \in \mathbb{R}} \mathbb{E}_X [(b - \sin(\pi X))^2].$$

Note that since $X \sim \text{Uniform}[-1, 1]$ we have that the probability density function for X is

$$f_X(x) = \begin{cases} \frac{1}{2} & x \in [-1, 1], \\ 0 & \text{otherwise,} \end{cases}$$

and thus

$$\begin{aligned}\mathbb{E}_X [(b - \sin(\pi X))^2] &= \int_{-1}^1 \frac{1}{2} (b - \sin(\pi x))^2 dx \\ &= b^2 + \frac{1}{2}.\end{aligned}$$

Taking a derivative with respect to b and setting this equal to zero yields the requirement that $2b^* = 0$, and hence $b^* = 0$, i.e., $h^\sharp(x) = 0$. We can also now easily compute

$$R(h^\sharp) = \mathbb{E}_X [(h^\sharp(X) - \sin(\pi X))^2] = \mathbb{E}_X [(\sin(\pi X))^2] = \frac{1}{2}.$$

The next step in our calculations is to determine $\bar{h}(X) = \mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(X)]$. To do this we must consider what $h_{\mathcal{D}}$ will be for any dataset $\mathcal{D} = \{(x_1, y_1), (x_2, y_2)\}$, where $y_i = \sin(\pi x_i)$ and $X_1, X_2 \sim \text{Uniform}[-1, 1]$. Given any particular \mathcal{D} , a natural strategy is to choose $h_{\mathcal{D}}(x)$ to minimize the empirical squared error, i.e., to choose b to minimize

$$(b - y_1)^2 + (b - y_2)^2.$$

It is straightforward to show that the optimal b is given by

$$b = \frac{y_1 + y_2}{2},$$

and thus

$$h_{\mathcal{D}}(X) = \frac{y_1 + y_2}{2} = \frac{\sin(\pi x_1) + \sin(\pi x_2)}{2}.$$

From this we have

$$\begin{aligned}\bar{h}(X) &= \mathbb{E}_{X_1, X_2} \left[\frac{\sin(\pi X_1) + \sin(\pi X_2)}{2} \right] \\ &= \mathbb{E}_{X_1} \left[\frac{\sin(\pi X_1)}{2} \right] + \mathbb{E}_{X_2} \left[\frac{\sin(\pi X_2)}{2} \right] \\ &= 0.\end{aligned}$$

The last equality can be seen either by symmetry or by manually calculating the expectation using an integral. Finally, we can conclude that

$$\begin{aligned}
 \text{bias}^2 &= \mathbb{E}_X [\text{bias}(h_{\mathcal{D}}(X))^2] \\
 &= \mathbb{E}_X [(\bar{h}(X) - h^*(X))^2] \\
 &= \mathbb{E}_X [(\sin(\pi X_2))^2] \\
 &= \frac{1}{2}.
 \end{aligned}$$

The final step is to compute the variance term. Specifically, we must compute

$$\mathbb{E}_X [\text{var}(h_{\mathcal{D}}(X)|X)] = \mathbb{E}_X [E_{\mathcal{D}} [(h_{\mathcal{D}}(X) - \bar{h}(X))^2]]$$

Note that

$$\begin{aligned}
 \mathbb{E}_{\mathcal{D}} [(h_{\mathcal{D}}(X) - \bar{h}(X))^2] &= \mathbb{E}_{X_1, X_2} \left[\left(\frac{\sin(\pi X_1) + \sin(\pi X_2)}{2} - 0 \right)^2 \right] \\
 &= \mathbb{E}_{X_1} \left[\frac{\sin^2(\pi X_1)}{4} \right] + \mathbb{E}_{X_2} \left[\frac{\sin^2(\pi X_2)}{4} \right] \\
 &\quad + \mathbb{E}_{X_1, X_2} [\sin(\pi X_1) \sin(\pi X_2)] \\
 &= \frac{1}{8} + \frac{1}{8} + \mathbb{E}_{X_1} [\sin(\pi X_1)] \mathbb{E}_{X_2} [\sin(\pi X_2)] \\
 &= \frac{1}{4}
 \end{aligned}$$

Thus

$$\mathbb{E}_X [\text{var}(h_{\mathcal{D}}(X)|X)] = \mathbb{E}_X \left[\frac{1}{4} \right] = \frac{1}{4}.$$

Putting this all together, we have

$$\mathbb{E}_{\mathcal{D}} [R(h_{\mathcal{D}})] = \frac{1}{2} + \frac{1}{4} = \frac{3}{4}.$$

Bias-variance decomposition for \mathcal{H}_1

As before, we will begin by computing

$$h^{\#} = \arg \min_{h \in \mathcal{H}_1} \mathbb{E}_X [(h - h^*(X))^2].$$

This is somewhat more complicated since now $h \in \mathcal{H}_1$ is an arbitrary affine function, thus we must solve

$$(a^*, b^*) = \arg \min_{a, b \in \mathbb{R}} \mathbb{E}_X [(aX + b - \sin(\pi X))^2].$$

Now, note that

$$\begin{aligned} \mathbb{E}_X [(aX + b - \sin(\pi X))^2] &= \int_{-1}^1 \frac{1}{2} (ax + b - \sin(\pi x))^2 dx \\ &= \frac{a^2}{3} - \frac{2a}{\pi} + b^2 + \frac{1}{2}. \end{aligned}$$

Taking a derivative with respect to b and setting it equal to zero again yields the requirement that $2b^* = 0$, and hence $b^* = 0$. Repeating with a yields the requirement that $\frac{2}{3}a^* - \frac{2}{\pi} = 0$, and hence $a^* = \frac{3}{\pi}$. Thus $h^{\#}(x) = \frac{3}{\pi}x$. We can also now easily compute

$$\begin{aligned} R(h^{\#}) &= \mathbb{E}_X [(h^{\#}(X) - \sin(\pi X))^2] \\ &= \mathbb{E}_X \left[\left(\frac{3}{\pi}X - \sin(\pi X) \right)^2 \right] \\ &= \frac{(3/\pi)^2}{3} - \frac{2(3/\pi)}{\pi} + \frac{1}{2} = \frac{1}{2} - \frac{3}{\pi^2} \approx 0.196. \end{aligned}$$

Note that this is substantially smaller than $R(h^\sharp)$ for \mathcal{H}_0 .

The next step in our calculations is to determine $\bar{h}(X) = \mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(X)]$. As before, we must consider what $h_{\mathcal{D}}$ will be for any dataset $\mathcal{D} = \{(x_1, y_1), (x_2, y_2)\}$, where $y_i = \sin(\pi x_i)$ and $X_1, X_2 \sim \text{Uniform}[-1, 1]$. If we follow the same approach of minimizing the empirical error, this will result in the a and b that *interpolate* the dataset. The formula for this will be given by

$$a = \frac{y_2 - y_1}{x_2 - x_1} \quad b = \frac{x_2 y_1 - x_1 y_2}{x_2 - x_1}.$$

The formula for a is simply the standard formula for the slope of a line passing through two points, and the formula for b can be obtained by simply constraining $y_1 = ax_1 + b$ for this choice of a . Thus

$$\begin{aligned} h_{\mathcal{D}}(X) &= \frac{y_2 - y_1}{x_2 - x_1} X + \frac{x_2 y_1 - x_1 y_2}{x_2 - x_1} \\ &= \frac{(\sin(\pi x_2) - \sin(\pi x_1))X + x_2 \sin(\pi x_1) - x_1 \sin(\pi x_2)}{x_2 - x_1}. \end{aligned}$$

From linearity we have

$$\begin{aligned} \bar{h}(X) &= \mathbb{E}_{X_1, X_2} [h_{\mathcal{D}}(X)] \\ &= \bar{a}X + \bar{b}, \end{aligned}$$

where

$$\begin{aligned} \bar{a} &= \mathbb{E}_{X_1, X_2} \left[\frac{(\sin(\pi X_2) - \sin(\pi X_1))}{X_2 - X_1} \right] \\ &= \int_{-1}^1 \int_{-1}^1 \frac{1}{4} \cdot \frac{(\sin(\pi x_2) - \sin(\pi x_1))}{x_2 - x_1} dx_1 dx_2 \\ &\approx 0.7759 \end{aligned}$$

and

$$\begin{aligned}\bar{b} &= \mathbb{E}_{X_1, X_2} \left[\frac{X_2 \sin(\pi X_1) - X_1 \sin(\pi X_2)}{X_2 - X_1} \right] \\ &= \int_{-1}^1 \int_{-1}^1 \frac{1}{4} \cdot \frac{(x_2 \sin(\pi x_1) - x_1 \sin(\pi x_2))}{x_2 - x_1} dx_1 dx_2 \\ &= 0.\end{aligned}$$

Thus,

$$\bar{h}(X) \approx 0.7759X.$$

Finally, we can conclude that

$$\begin{aligned}\text{bias}^2 &= \mathbb{E}_X [\text{bias}(h_{\mathcal{D}}(X))^2] \\ &= \mathbb{E}_X [(\bar{h}(X) - h^*(X))^2] \\ &= \mathbb{E}_X [(0.7759X - \sin(\pi X))^2] \\ &= \frac{(0.7759)^2}{3} - \frac{2 \cdot 0.7759}{\pi} + \frac{1}{2} \approx 0.207\end{aligned}$$

The final step is to compute the variance term. Specifically, we must compute

$$\mathbb{E}_X [\text{var}(h_{\mathcal{D}}(X)|X)] = \mathbb{E}_X [E_{\mathcal{D}} [(h_{\mathcal{D}}(X) - \bar{h}(X))^2]]$$

Note that if we let

$$A_{\mathcal{D}} = \frac{\sin(\pi X_2) - \sin(\pi X_1)}{X_2 - X_1}$$

and

$$B_{\mathcal{D}} = \frac{X_2 \sin(\pi X_1) - X_1 \sin(\pi X_2)}{X_2 - X_1},$$

then we can write

$$\begin{aligned}\mathbb{E}_{\mathcal{D}} [(h_{\mathcal{D}}(X) - \bar{h}(X))^2] &= \mathbb{E}_{\mathcal{D}} [(A_{\mathcal{D}}X + B_{\mathcal{D}} - 0.7759X)^2] \\ &= X^2 \mathbb{E}_{\mathcal{D}} [(A_{\mathcal{D}} - 0.7759)^2] \\ &\quad + 2X \mathbb{E}_{\mathcal{D}} [(A_{\mathcal{D}} - 0.7759)B_{\mathcal{D}}] + \mathbb{E}_{\mathcal{D}} [B_{\mathcal{D}}^2].\end{aligned}$$

Note that

$$\begin{aligned}\mathbb{E}_{\mathcal{D}} [(A_{\mathcal{D}} - 0.7759)^2] &= (\mathbb{E}_{\mathcal{D}} [A_{\mathcal{D}}^2] - 0.7759 \cdot \mathbb{E}_{\mathcal{D}} [A_{\mathcal{D}}] + 0.7759^2) \\ &\approx (\mathbb{E}_{\mathcal{D}} [A_{\mathcal{D}}^2] - 0.7759^2),\end{aligned}$$

where the final (approximate) equality follows from the fact that $\mathbb{E}_{\mathcal{D}} [A_{\mathcal{D}}] = \bar{a} \approx 0.7759$. We can compute

$$\begin{aligned}\mathbb{E}_{\mathcal{D}} [A_{\mathcal{D}}^2] &= \int_{-1}^1 \int_{-1}^1 \frac{1}{4} \cdot \left(\frac{(\sin(\pi x_2) - \sin(\pi x_1))}{x_2 - x_1} \right)^2 dx_1 dx_2 \\ &\approx 2.8981.\end{aligned}$$

Similarly, we can also compute

$$\begin{aligned}\mathbb{E}_{\mathcal{D}} [B_{\mathcal{D}}^2] &= \int_{-1}^1 \int_{-1}^1 \frac{1}{4} \cdot \left(\frac{(x_2 \sin(\pi x_1) - x_1 \sin(\pi x_2))}{x_2 - x_1} \right)^2 dx_1 dx_2 \\ &\approx 0.9109.\end{aligned}$$

We will not bother to compute

$$\rho = \mathbb{E}_{\mathcal{D}} [(A_{\mathcal{D}} - 0.7759)B_{\mathcal{D}}]$$

since, as we will see below, it does not matter. Putting this together yields

$$\begin{aligned}\mathbb{E}_X [\text{var}(h_{\mathcal{D}}(X)|X)] &\approx \mathbb{E}_X [(2.8981 - 0.7759^2)X^2 + 2\rho X + 0.9109] \\ &\approx 2.2961 \mathbb{E}_X [X^2] + 2\rho \mathbb{E}_X [X] + 0.9109 \\ &= \frac{2.2961}{3} + 0.9109 \approx 1.676,\end{aligned}$$

where above we have used that for $X \sim \text{Uniform}[-1, 1]$, $\mathbb{E}[X] = 0$ and $\mathbb{E}[X^2] = \frac{1}{3}$. Putting this all together, we have

$$\mathbb{E}_{\mathcal{D}} [R(h_{\mathcal{D}})] \approx 0.207 + 1.676 = 1.883.$$

Thus, despite the fact that \mathcal{H}_1 resulted in a much smaller $R(h^\#)$, we see that $\mathbb{E}_{\mathcal{D}} [R(h_{\mathcal{D}})]$ is *substantially* larger than before because of the much larger variance.