

The bias variance tradeoff

The world we think of

Data-point :	$\{x, y\} \sim p(x, y)$
Data-set :	$D = (X, Y) \sim p^n$
Model class :	\mathcal{A}
Model :	$h_D(x) \in \mathcal{A}$
Expected model :	$\bar{h}(x) = \int_D h_D \Pr(D) \, dh_D = E_D[h_D]$
Expected label :	$\bar{y}(x) = \int_y y \Pr(y x) \, dy = E_{y x}[y]$

The bias variance tradeoff is then encapsulated in

$$\underbrace{E_{x,y,D} \left[(h_D(x) - y)^2 \right]}_{\text{Squared error}} = \underbrace{E_{x,D} \left[(h_D(x) - \bar{h}(x))^2 \right]}_{\text{variance}} \\
 + \underbrace{E_x \left[(\bar{h}(x) - \bar{y}(x))^2 \right]}_{\text{bias}^2} \\
 + \underbrace{E_{x,y} \left[(\bar{y}(x) - y)^2 \right]}_{\text{noise}}$$

This is all explained in much more depth by the creators of CS4780 at Cornell
<http://www.cs.cornell.edu/courses/cs4780/2017sp/lectures/lecturenote11.html>

So where's the tradeoff?

If we have $y = f(x) + e$ and consider measure related to the total error, the MSE, defined as

$$MSE(x) = \text{squared error} - \text{noise}$$

We have a nice relationship

$$MSE = \text{Var} + \text{bias}^2 \quad \Rightarrow$$

$$\sqrt{MSE}^2 = \sigma^2 + \text{bias}^2 \quad \Rightarrow$$

$$RMSE^2 = \sigma^2 + \text{bias}^2$$

It turns out that σ and bias are orthogonal vectors \Rightarrow this is a right-angled triangle.

Imagine $RMSE$ is minimized (no low hanging fruits to independently lower either bias or variance). The relation between bias and σ is then similar to how two legs in a triangle can be scaled while preserving the length of the hypotenuse; clearly a tradeoff.