

The bias-variance tradeoff is an important aspect of ML
 Consider the additive error model $y = f(x) + \epsilon$ where $\epsilon \sim N(0, \sigma^2)$. Here we are thinking about choosing an algorithm to learn the predictive $f(\cdot)$.
 All learning algorithms use a mathematical approach that contains an 'error' term which can be further split into two components:

$$\text{EPE} = \mathbb{E}(\epsilon_0 - \hat{f}(x_0))^2 = \underbrace{\sigma^2}_{\text{irreducible error}} + \mathbb{E}(f(x_0) - \hat{f}(x_0))^2$$

$\epsilon_0 = f(x_0) + \epsilon_0 = \sigma^2 + \text{Bias}^2 + \text{Variance}$

This is the noise term in the true relationship that cannot fundamentally be reduced by the model error
Irreducible error
Reducible error
should be minimized further to max. accuracy

$$\text{Claim: } \mathbb{E}(\epsilon_0 - \hat{f}(x_0))^2 = \sigma^2 + \text{Bias}^2(\hat{f}(x_0)) + \text{Var}(\hat{f}(x_0))$$

Decompose the expected prediction error into 3 parts: irreducible error, squared Bias, variance

Proof:
$$\text{EPE} = \mathbb{E}\{(y_0 - \hat{f}(x_0))^2\}$$

$$= \mathbb{E}\{(f(x_0) + \epsilon_0 - \hat{f}(x_0))^2\}$$

$$= \mathbb{E}\{\epsilon_0^2 + 2\epsilon_0(f(x_0) - \hat{f}(x_0)) + (f(x_0) - \hat{f}(x_0))^2\}$$

$$= \mathbb{E}\epsilon_0^2 + 2\mathbb{E}\{\epsilon_0(f(x_0) - \hat{f}(x_0))\} + \mathbb{E}\{(f(x_0) - \hat{f}(x_0))^2\}$$

$$\mathbb{E}\epsilon_0^2 = \sigma^2$$

$$\mathbb{E}(xy) = 0 \text{ if } x \perp y \text{ as } \mathbb{E}x = 0$$

$$= \sigma^2 + 0 + \mathbb{E}\{(f(x_0) - \hat{f}(x_0))^2\}$$

$$\mathbb{E}\{(f(x_0) - \hat{f}(x_0))^2\} = \mathbb{E}\{(f(x_0) - \mathbb{E}\hat{f}(x_0) + \mathbb{E}\hat{f}(x_0) - \hat{f}(x_0))^2\}$$

$$= \mathbb{E}\{(f(x_0) - \mathbb{E}\hat{f}(x_0))^2 + 2(f(x_0) - \mathbb{E}\hat{f}(x_0))(\mathbb{E}\hat{f}(x_0) - \hat{f}(x_0)) + (\mathbb{E}\hat{f}(x_0) - \hat{f}(x_0))^2\}$$

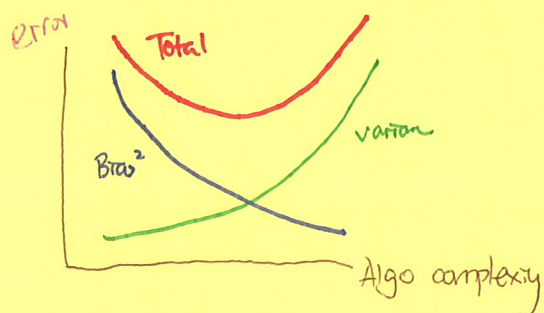
$$= (f(x_0) - \mathbb{E}\hat{f}(x_0))^2 + 2\mathbb{E}(\dots) + \mathbb{E}(\mathbb{E}\hat{f}(x_0) - \hat{f}(x_0))^2$$

$$\mathbb{E}\{(f(x_0) - \hat{f}(x_0))^2\} = (\mathbb{E}\hat{f}(x_0) - f(x_0))^2 + \mathbb{E}\{(\hat{f}(x_0) - \mathbb{E}\hat{f}(x_0))^2\}$$

$$= \text{Bias}^2(\hat{f}(x_0)) + \text{Var}(\hat{f}(x_0))$$

error due to overgeneralisation data from training set
error due to fluctuation in the training dataset from test set

this is often called Bias-variance tradeoff because improving one will worsen the other



In ML, the ideal algorithm has low bias and can accurately model the true relationship and has low variability, by producing consistent prediction across different datasets.