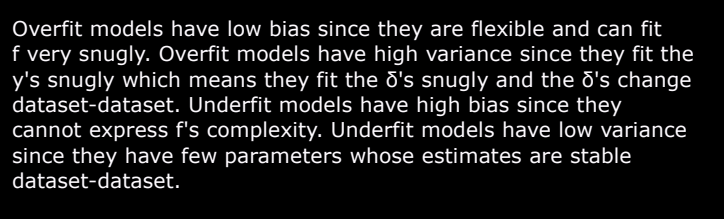

$$MSE(\alpha) = E_{\Delta_1, \Delta_2, \dots, \Delta_h, \Delta_{\alpha}} [(Y_{\alpha} - g(\hat{x}_{\alpha}))^2 | X, \hat{x}_{\alpha}]$$

Let's step it up two more notches in one shot (a) Let x_1, \dots, x_n be random realizations from $P(X)$, thus X is random and (b) x^* is a random realization from $P(X)$ as well.

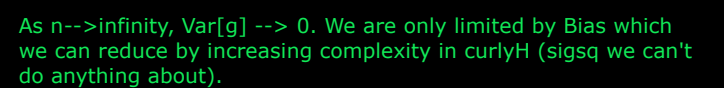
This is the general bias-variance decomposition formula. This is not computable in practice since you only have one D and one g . It is a theoretical formula.

Is there a "bias-variance tradeoff"? Yes and no...
It is not a zero sum game. If you make smart modeling decisions both bias AND variance decrease.



MSE look like over n ?

The graph shows two curves representing Mean Squared Error (MSE) as a function of sample size n . The y-axis is labeled 'MSE' and the x-axis is labeled 'n'. A green curve, labeled 'OOS' (Out-of-Sample), starts high and decreases as n increases, approaching a horizontal dashed line. A red curve, labeled 'in-sample', starts lower and increases as n increases, also approaching the same horizontal dashed line. The vertical distance between the two curves is labeled 'Bias²' with an arrow. The horizontal dashed line represents the true parameter value.



Let's consider a metaalgorithm called "model averaging". You fit g_1, g_2, \dots, g_M and then you ship their average,

$$g_{g \times g} := \frac{g_1 + g_2 + \dots + g_m}{m}$$

What is the MSE of g_avg ?

$$MSE = \sigma^2 + E_x \left[\text{Bias}[\hat{\theta}_{avg}]^2 \right] + E_x \left[\text{Var}[\hat{\theta}_{avg}] \right]$$

$$-x \left[\begin{pmatrix} m_1 & 0 & \dots & 0 \\ 0 & m_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & m_m \end{pmatrix} \right], \quad m^2 \left[x \left[\sqrt{g_1 + \dots + g_m} \right] \right]$$

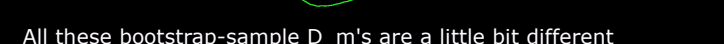
assume (b) g_1, \dots, g_M are independent models and (c)
The variance of each are the same too.

minimum

dependence on one another. (a) can be made true in practice since we can just use algorithms that overfit. (c) doesn't matter if (b) is true since that term $\rightarrow 0$ as $M \rightarrow \infty$. The problem is (b)!!!! And it seems to be impossible to solve!!

Consider this: In 1964, 60% of the U.S. population (that's

Math 241) that approximately $2/3$ of the rows of D appear in each sample and $1/3$ are missing. Do this M times.



and different observations duplicated. Then, you fit your models using the same A and H,

$$g_1 = A(\mathbb{D}_1, \mathcal{A}), g_2 = A(\mathbb{D}_2, \mathcal{A}), \dots, g_m = A(\mathbb{D}_m, \mathcal{A})$$

you didn't use bootstrap samples of D . So when you average,

$$0 \leq m \leq \sum_{b=1}^n v_b$$

the reduction from bagging in MSE next class.