

Bias-Variance trade-off

DRIPTA MJ

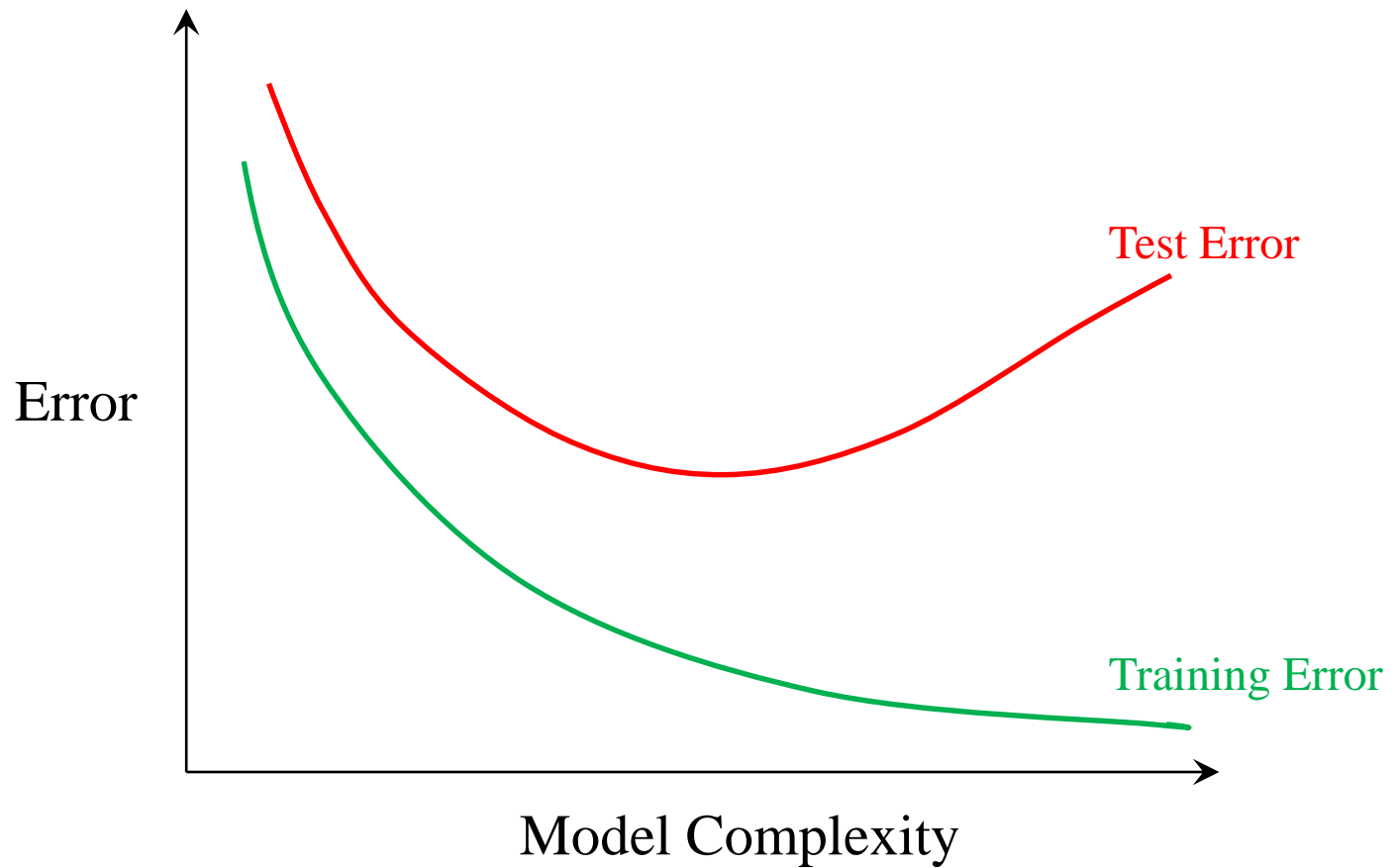
Department of Mathematics

RAMAKRISHNA MISSION VIVEKANANDA EDUCATIONAL AND RESEARCH INSTITUTE
BELUR MATH, INDIA

Machine Learning
CS230

Sem 3, 2018-19

Error vs complexity



Bias-variance decomposition

- Dataset: $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$
- Let $g_{\mathcal{D}}$ be the hypothesis which is fit to a particular training dataset \mathcal{D}
- Want to compute the expected prediction error at an arbitrary test point with input \mathbf{x} and output y : $\mathbb{E}_{\mathbf{x}, y, \mathcal{D}} \left[(g_{\mathcal{D}}(\mathbf{x}) - y)^2 \right]$.

- Mean prediction of the machine learning algorithm:

$$\bar{g}(\mathbf{x}) = \mathbb{E}_{\mathcal{D}} \left[g_{\mathcal{D}}(\mathbf{x}) \right]$$

- So determining the value of $\bar{g}(\mathbf{x})$ involve
 - generating different training datasets (\mathcal{D}),
 - training separate functions ($g_{\mathcal{D}}$) for every generated dataset,
 - making predictions at an arbitrary test point \mathbf{x} with all trained functions,
 - and finally, averaging over all the predictions.
- Let $\bar{y}(\mathbf{x})$ be the expected value of the output at \mathbf{x} , i.e. $\bar{y}(\mathbf{x}) = \mathbb{E}_{y|\mathbf{x}}[y]$.

Bias-variance decomposition

- The expected error can be simplified as:

$$\begin{aligned}\mathbb{E}_{\mathbf{x},y,\mathcal{D}} \left[(g_{\mathcal{D}}(\mathbf{x}) - y)^2 \right] &= \mathbb{E}_{\mathbf{x},y,\mathcal{D}} \left[\left((g_{\mathcal{D}}(\mathbf{x}) - \bar{g}(\mathbf{x})) + (\bar{g}(\mathbf{x}) - y) \right)^2 \right] \\ &= \mathbb{E}_{\mathbf{x},\mathcal{D}} \left[(g_{\mathcal{D}}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2 \right] + 2\mathbb{E}_{\mathbf{x},y,\mathcal{D}} \left[(g_{\mathcal{D}}(\mathbf{x}) - \bar{g}(\mathbf{x})) (\bar{g}(\mathbf{x}) - y) \right] \\ &\quad + \mathbb{E}_{\mathbf{x},y} \left[(\bar{g}(\mathbf{x}) - y)^2 \right]\end{aligned}$$

The second term on simplification yields

$$\begin{aligned}2\mathbb{E}_{\mathbf{x},y,\mathcal{D}} \left[(g_{\mathcal{D}}(\mathbf{x}) - \bar{g}(\mathbf{x})) (\bar{g}(\mathbf{x}) - y) \right] &= 2\mathbb{E}_{\mathbf{x},y} \left[\mathbb{E}_{\mathcal{D}} \left[(g_{\mathcal{D}}(\mathbf{x}) - \bar{g}(\mathbf{x})) \right] (\bar{g}(\mathbf{x}) - y) \right] \\ &= 2\mathbb{E}_{\mathbf{x},y} \left[\left(\mathbb{E}_{\mathcal{D}} \left[g_{\mathcal{D}}(\mathbf{x}) \right] - \bar{g}(\mathbf{x}) \right) (\bar{g}(\mathbf{x}) - y) \right] \\ &= 2\mathbb{E}_{\mathbf{x},y} \left[(\bar{g}(\mathbf{x}) - \bar{g}(\mathbf{x})) (\bar{g}(\mathbf{x}) - y) \right] \\ &= 2\mathbb{E}_{\mathbf{x},y} [0] \\ &= 0\end{aligned}$$

Bias-variance decomposition

The third term can be simplified as

$$\begin{aligned}\mathbb{E}_{\mathbf{x},y} \left[(\bar{g}(\mathbf{x}) - y)^2 \right] &= \mathbb{E}_{\mathbf{x}} \left[(\bar{g}(\mathbf{x}) - \bar{y}(\mathbf{x}))^2 \right] + \mathbb{E}_{\mathbf{x},y} \left[(\bar{y}(\mathbf{x}) - y)^2 \right] \\ &\quad + 2\mathbb{E}_{\mathbf{x},y} \left[(\bar{g}(\mathbf{x}) - \bar{y}(\mathbf{x})) (\bar{y}(\mathbf{x}) - y) \right]\end{aligned}$$

where $\bar{y}(\mathbf{x}) = \mathbb{E}_{y|\mathbf{x}}[y]$ and the last term can be simplified as:

$$\begin{aligned}2\mathbb{E}_{\mathbf{x},y} \left[(\bar{g}(\mathbf{x}) - \bar{y}(\mathbf{x})) (\bar{y}(\mathbf{x}) - y) \right] &= 2\mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{y|\mathbf{x}}[(\bar{y}(\mathbf{x}) - y)] (\bar{g}(\mathbf{x}) - \bar{y}(\mathbf{x})) \right] \\ &= 2\mathbb{E}_{\mathbf{x}} \left[(\bar{y}(\mathbf{x}) - \mathbb{E}_{y|\mathbf{x}}[y]) (\bar{g}(\mathbf{x}) - \bar{y}(\mathbf{x})) \right] \\ &= 2\mathbb{E}_{\mathbf{x}} \left[(\bar{y}(\mathbf{x}) - \bar{y}(\mathbf{x})) (\bar{g}(\mathbf{x}) - \bar{y}(\mathbf{x})) \right] \\ &= 0\end{aligned}$$

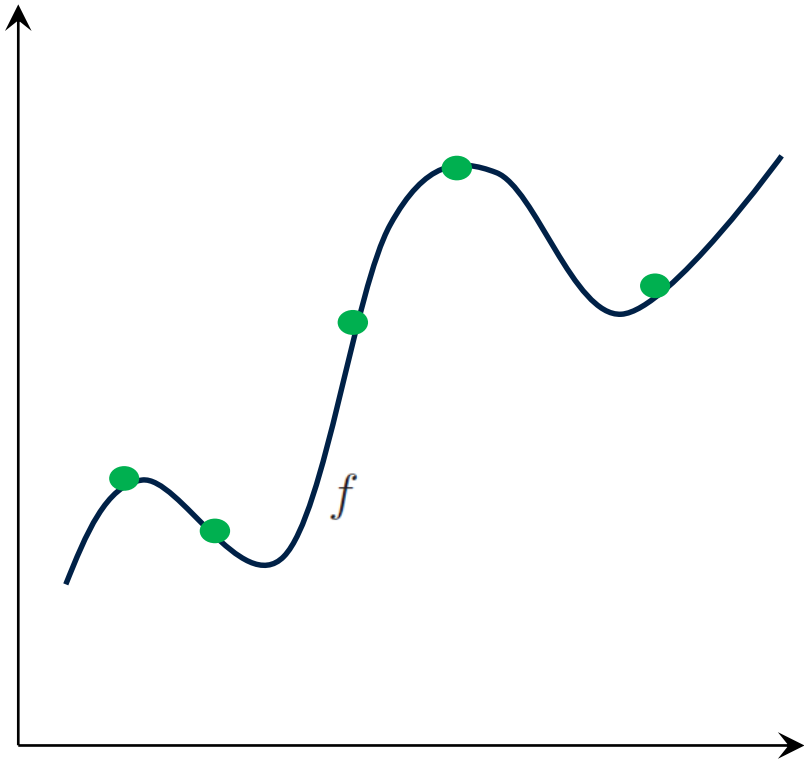
Bias-variance decomposition

$$\mathbb{E}_{\mathbf{x}, y, \mathcal{D}} \left[(g_{\mathcal{D}}(\mathbf{x}) - y)^2 \right] = \underbrace{\mathbb{E}_{\mathbf{x}, \mathcal{D}} \left[(g_{\mathcal{D}}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2 \right]}_{\text{Variance}} + \underbrace{\mathbb{E}_{\mathbf{x}} \left[(\bar{g}(\mathbf{x}) - \bar{y}(\mathbf{x}))^2 \right]}_{\text{Bias}^2} + \underbrace{\mathbb{E}_{\mathbf{x}, y} \left[(\bar{y}(\mathbf{x}) - y)^2 \right]}_{\text{Noise}}$$

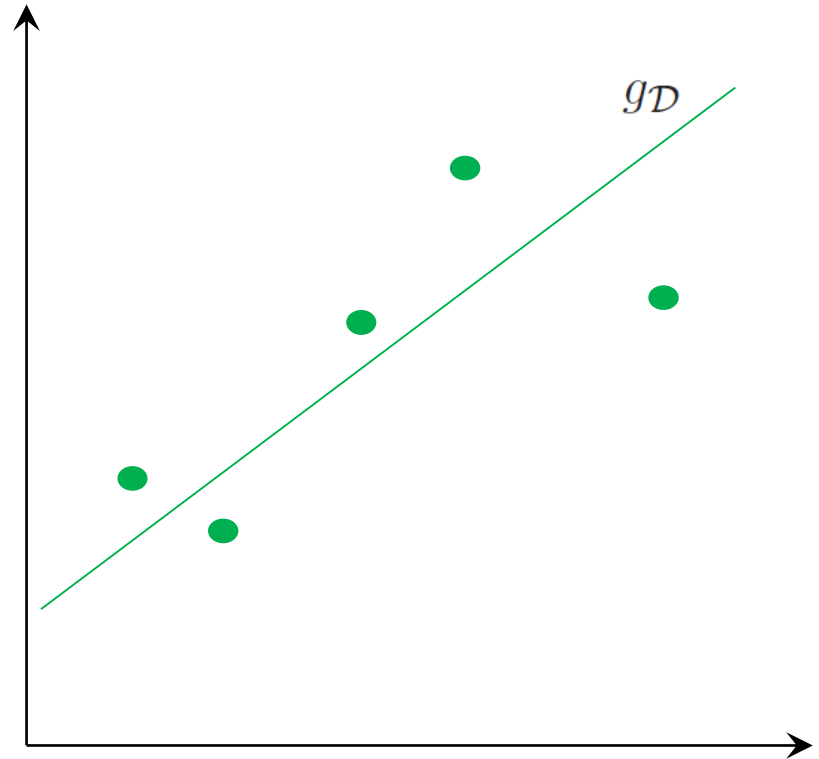
- **Variance**: It expresses the sensitivity of the solution on the particular choice of dataset \mathcal{D} .
- **Bias**: Difference between the expected prediction (averaged over different datasets) and the expected output value. This is the inherent error arising from the choice of model.
- **Noise**: Expresses the noise in the data.

Example

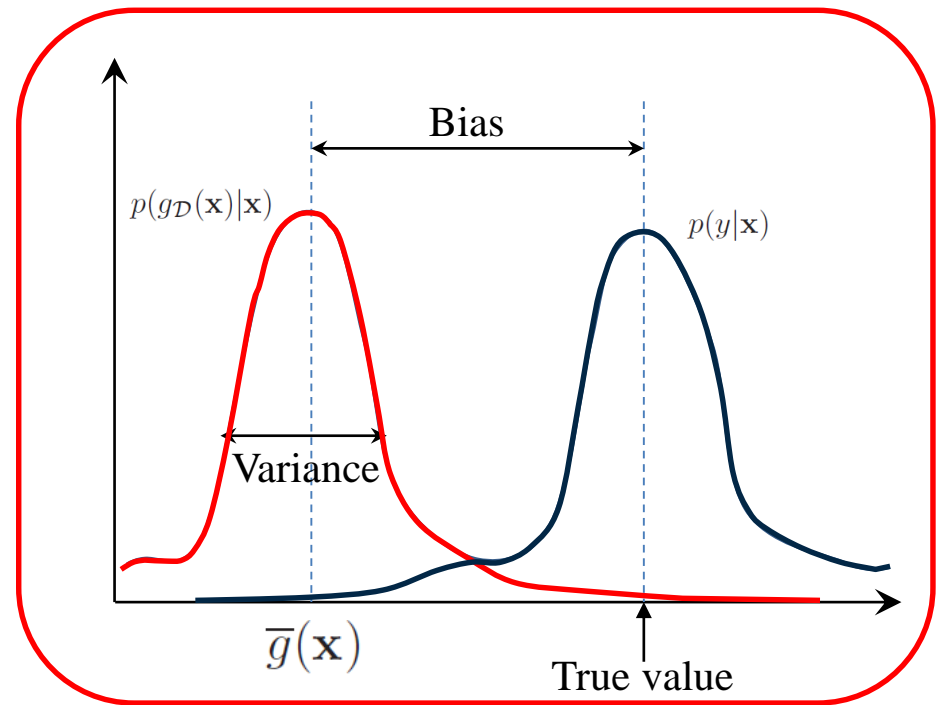
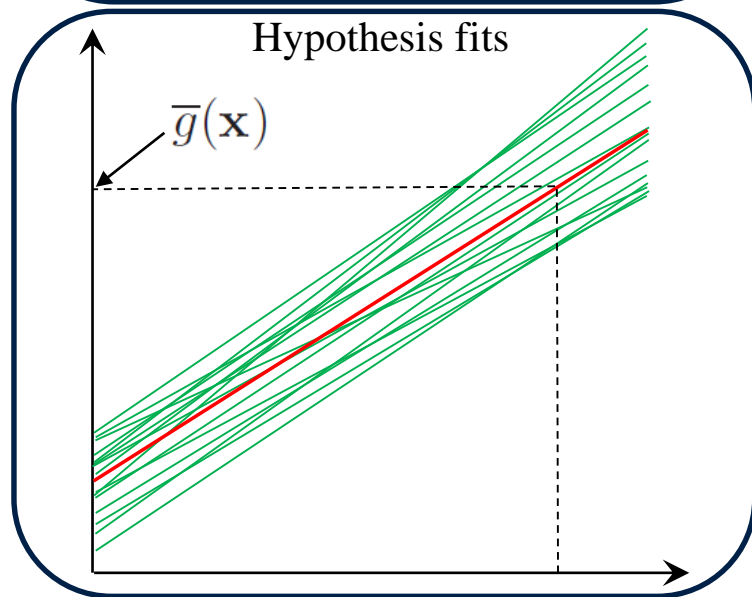
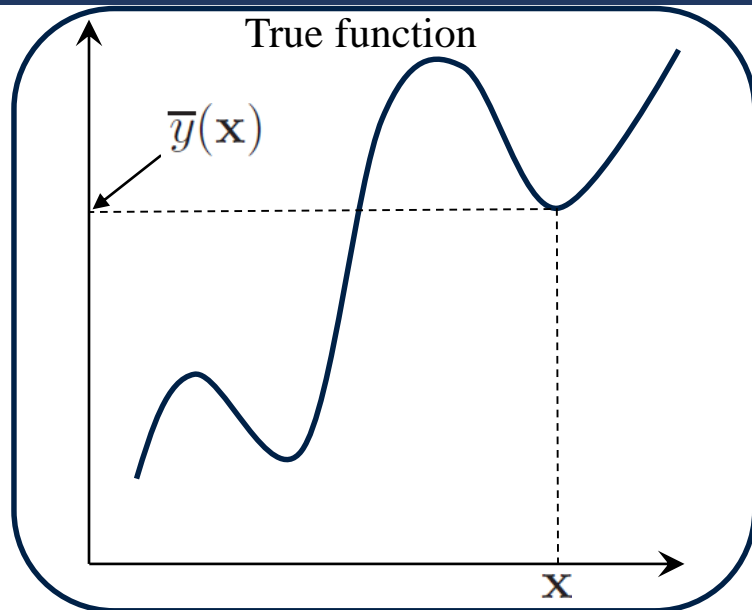
Underlying true function f
and data points



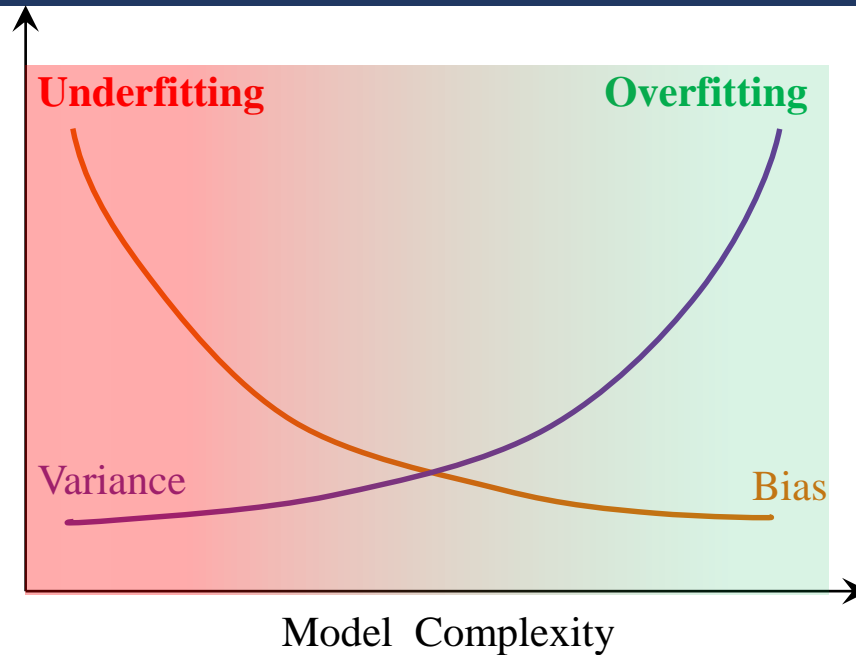
Hypothesis fit: $g_{\mathcal{D}}$



Visualization

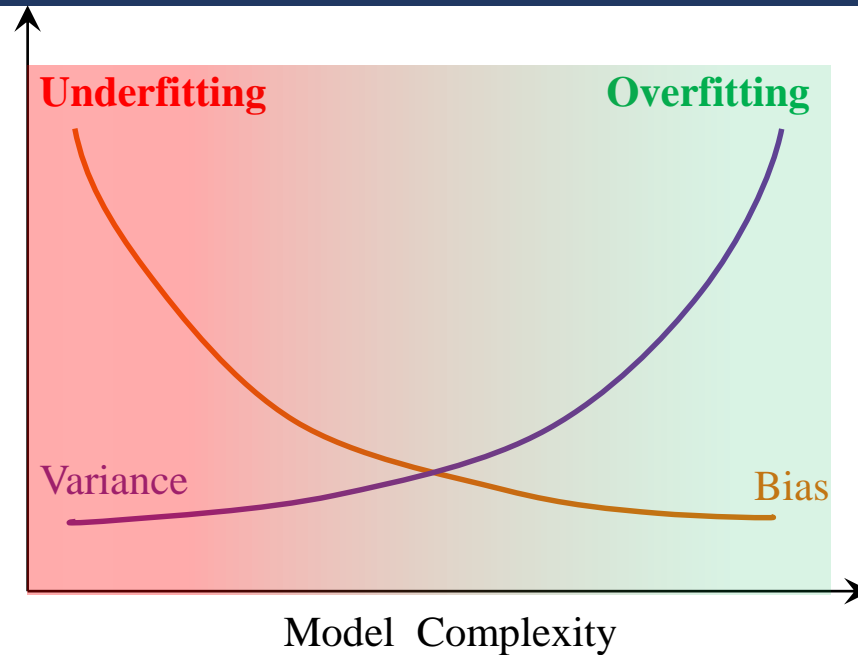


Bias, variance vs model complexity



- **High Bias:** Model is too simple, and so unable to fit the data properly.
 - Results in underfitting.
 - Training and test errors are both large.
- **High Variance:** Model is too complex, and so small changes in the data produce significant changes in the solution.
 - Results in overfitting.
 - Test Error \gg Training Error

Underfitting & Overfitting



- **Underfitting** can be addressed by
 - Increasing the complexity of the model.
 - Minimizing the cost function properly in the training stage.
- **Overfitting** can be addressed by
 - Reducing the complexity of the model.
 - Incorporating some form of regularization inside the cost function.