

Bias - Variance Tradeoff in ML

→ Used to understand the tradeoff between the simplicity of the model and its generalizability.

Bias: Error in model predictions as ML approximates a real-world variable with a simple, presupposed model architecture.

Simpler models introduce higher bias since typically, they have fewer parameters

Variance: Error in prediction due to a model's sensitivity

small fluctuations in the training dataset, since a model may be too complex and thus captures noise as a signal.

Bias

Leads to underfitting as the model misses important relationships between the target variable & the features.

Variance

Leads to overfitting as the model performs poorly on test data due to an overly complex architecture

Trade-off: To find an optimally complex model

$$\rightarrow (\text{Total error}) = (\text{Bias})^2 + (\text{Variance}) + (\text{Irreducible error})$$

Let y denote the real-world measurements of the target variable
 f denote the true model
 ϵ error in the true model \rightarrow assumed to be normally distributed.

$$y = f + \epsilon$$

$$\epsilon \sim N(0, \sigma^2)$$

\rightarrow estimate of irreducible error that even a true model will incur

Total error



Squared error

$$E[(y - \hat{f})^2]$$

best-fit model using a
given dataset.

Error

$$= E[(y - \hat{f})^2]$$

$$= E[y^2 + \hat{f}^2 - 2y\hat{f}]$$

$$= E[y^2] + E[\hat{f}^2] - E[2y\hat{f}]$$

$$\text{var}(\hat{f}) = E[(\hat{f} - E(\hat{f}))^2]$$

$$= E[(E(\hat{f}))^2 + \hat{f}^2 - 2\hat{f}E(\hat{f})]$$

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

$$\text{Var}(\hat{f}) = (E(\hat{f}))^2 + E(\hat{f}^2) - 2(E(\hat{f}))^2 = E(\hat{f}^2) - (E(\hat{f}))^2$$

$$\longrightarrow E(\hat{f}^2) = \text{Var}(\hat{f}) + (E(\hat{f}))^2$$

$$\text{Squared Error} = E[y^2] + \text{Var}(\hat{f}) + (E(\hat{f}))^2 - E[2y\hat{f}]$$

We need $E[y^2]$. Recall that $y = f + \epsilon$ → error

$$E[y^2] = E[(f + \epsilon)^2] = E[f^2] + E[2f\epsilon] + E[\epsilon^2]$$

Since the analysis is done for a given data point. Thus,

$$E[y^2] = f^2 + \underline{E(e) \cdot 2f} + E[e^2]$$

Recall that $e \sim N(0, \sigma^2)$

$$E(e) = 0$$

$$\sigma^2 = E[(e - \cancel{E(e)})^2] = E(e^2)$$

$$\boxed{\cancel{E[y^2]} = f^2 + 0 + \sigma^2 = (f^2 + \sigma^2)}$$

$$\begin{aligned} E[2y\hat{f}] &= 2E[(f+e) \cdot \hat{f}] = 2\left[E(f\hat{f}) + \underline{E(e\hat{f})}\right] \\ &= 2\left[fE(\hat{f}) + \cancel{E(e)}E(\hat{f})\right] \end{aligned}$$

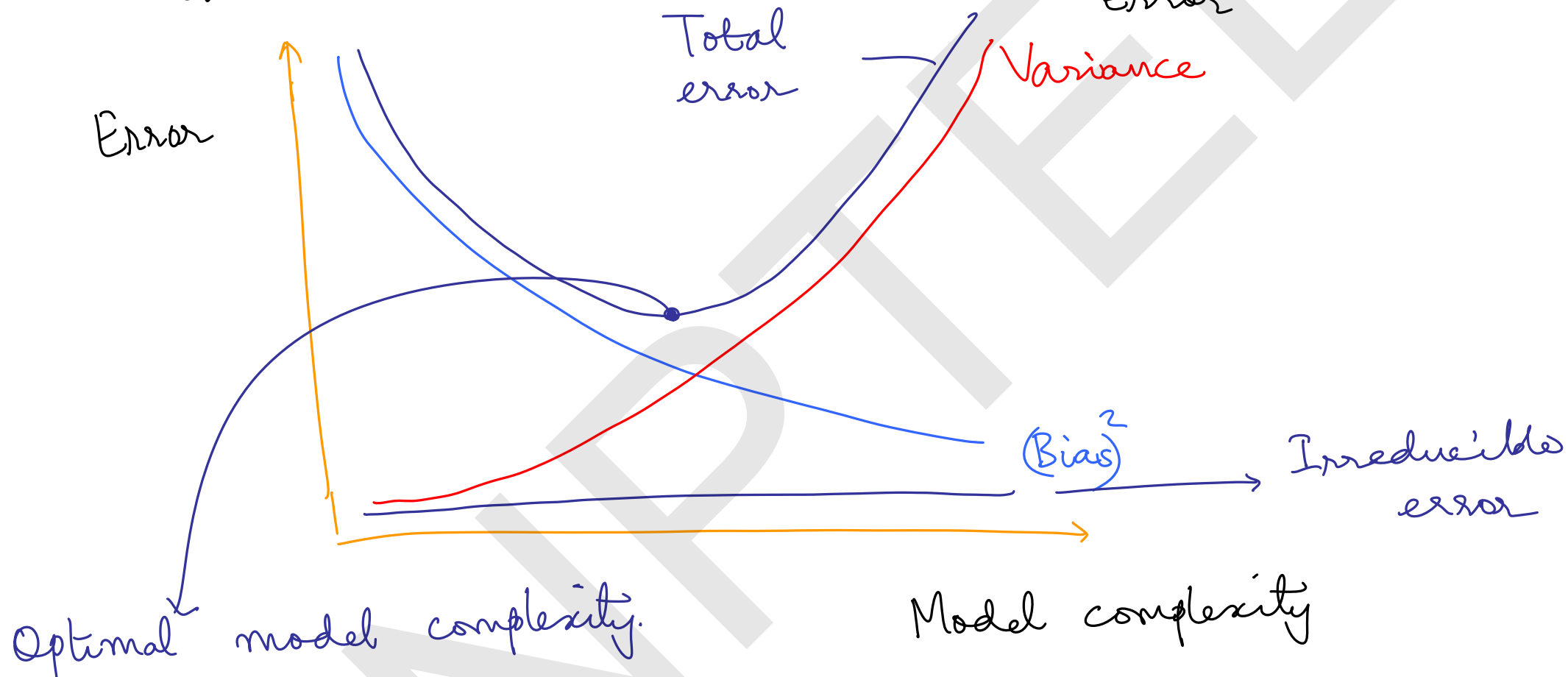
error for a datapoint and model prediction are independent

$$E[2y\hat{f}] = 2f E(\hat{f})$$

$$\begin{aligned} \text{Squared error} &= E[y^2] + \text{Var}(\hat{f}) + (E(\hat{f}))^2 - 2E[y\hat{f}] \\ &= (f^2 + \sigma^2) + \text{Var}(\hat{f}) + (E(\hat{f}))^2 - 2f E(\hat{f}) \\ &= \underbrace{(E[f - \hat{f}])^2}_{\text{Bias}^2} + \text{Var}(\hat{f}) + \sigma^2 \end{aligned}$$

$$\text{Squared error} = (\text{Bias}(\hat{f}))^2 + \text{Var}(\hat{f}) + \sigma^2$$

$$\text{Total error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$



To decrease bias: Increase model complexity

To decrease variance: Employing dimensionality reduction,
regularization, or feature selection