

Machine Learning

CMPT 726

Mo Chen

SFU School of Computing Science

2022-10-13

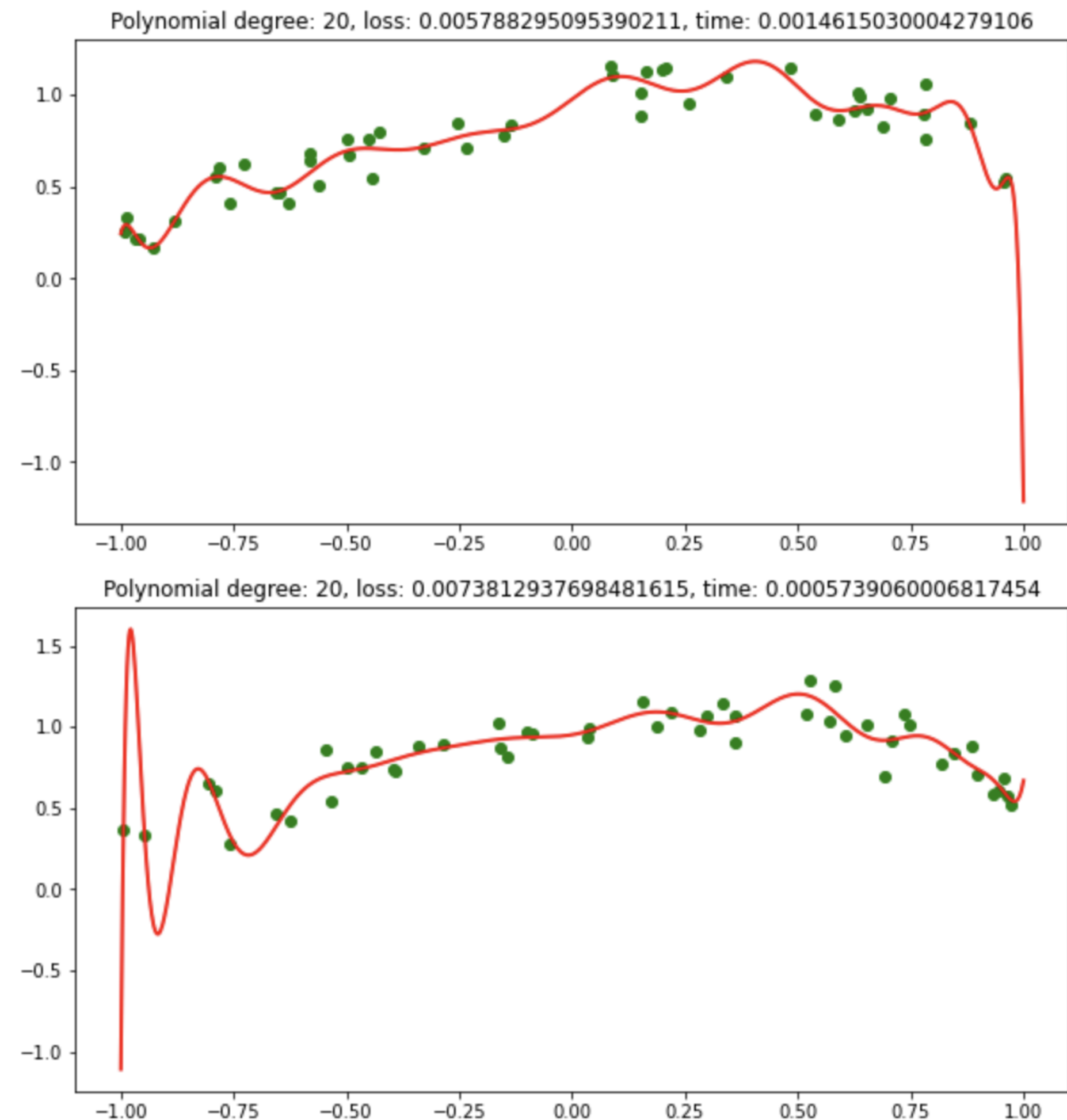
Bias-Variance Tradeoff

Understanding Overfitting

Recall: Overfitting happens when the model is fitting to the noise in the training data.

So, when overfitting happens, different training datasets can result in very different optimal parameter values, which often result in wrong predictions for some inputs. (See previous [Colab notebook](#))

Ideally, we would like to use a family of models that would typically produce accurate predictions on unseen testing data, regardless of the particular training dataset.



Bias-Variance Decomposition

$$\begin{aligned}\epsilon(\vec{x}, f, L) &= E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}] && \text{Expected Error} \\ &= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2 && \text{Bias Squared} \\ &\quad + \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) && \text{Variance} \\ &\quad + \text{Var}(y | \vec{x}) && \text{Irreducible Error}\end{aligned}$$

Bias-Variance Decomposition

$$\underbrace{\epsilon(\vec{x}, f, L)}_{\text{Unseen Testing Example}} = E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}] \quad \text{Expected Error}$$

$$= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2 \quad \text{Bias Squared}$$

$$+ \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) \quad \text{Variance}$$

$$+ \text{Var}(y | \vec{x}) \quad \text{Irreducible Error}$$

Bias-Variance Decomposition

$$\begin{aligned}\epsilon(\underbrace{\vec{x}, f, L}_{\text{Model}}) &= E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}] && \text{Expected Error} \\ &= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2 && \text{Bias Squared} \\ &\quad + \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) && \text{Variance} \\ &\quad + \text{Var}(y | \vec{x}) && \text{Irreducible Error}\end{aligned}$$

Bias-Variance Decomposition

$$\begin{aligned} \underbrace{\epsilon(\vec{x}, f, L)}_{\text{Loss Function}} &= E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}] && \text{Expected Error} \\ &= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2 && \text{Bias Squared} \\ &\quad + \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) && \text{Variance} \\ &\quad + \text{Var}(y | \vec{x}) && \text{Irreducible Error} \end{aligned}$$

Bias-Variance Decomposition

$$\begin{aligned}\epsilon(\vec{x}, f, L) &= E_{y, \mathcal{D}}[(f(\vec{x}; \underbrace{\theta^*(\mathcal{D}, L)}_{\text{Training Dataset}}) - y)^2 | \vec{x}] && \text{Expected Error} \\ &= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2 && \text{Bias Squared} \\ &\quad + \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) && \text{Variance} \\ &\quad + \text{Var}(y | \vec{x}) && \text{Irreducible Error}\end{aligned}$$

Bias-Variance Decomposition

$$\begin{aligned}\epsilon(\vec{x}, f, L) &= E_{y, \mathcal{D}}[(f(\vec{x}; \underbrace{\theta^*(\mathcal{D}, L)}_{\text{Optimal parameters on training data}}) - y)^2 | \vec{x}] && \text{Expected Error} \\ &= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2 && \text{Bias Squared} \\ &\quad + \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) && \text{Variance} \\ &\quad + \text{Var}(y | \vec{x}) && \text{Irreducible Error}\end{aligned}$$

Bias-Variance Decomposition

$$\begin{aligned}\epsilon(\vec{x}, f, L) &= E_{y, \mathcal{D}}[\underbrace{(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2}_{\text{Prediction of trained model on new testing example}} | \vec{x}] && \text{Expected Error} \\ &= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2 && \text{Bias Squared} \\ &\quad + \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) && \text{Variance} \\ &\quad + \text{Var}(y | \vec{x}) && \text{Irreducible Error}\end{aligned}$$

Bias-Variance Decomposition

$$\begin{aligned}\epsilon(\vec{x}, f, L) &= E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - \underbrace{y}_{\text{(Possibly noisy) label corresponding to testing example}})^2 | \vec{x}] && \text{Expected Error} \\ &= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2 && \text{Bias Squared} \\ &\quad + \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) && \text{Variance} \\ &\quad + \text{Var}(y | \vec{x}) && \text{Irreducible Error}\end{aligned}$$

Bias-Variance Decomposition

$$\begin{aligned}\epsilon(\vec{x}, f, L) &= E_{y, \mathcal{D}}[\underbrace{(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2}_{\text{Mean squared error (MSE) on testing example}} | \vec{x}] && \text{Expected Error} \\ &= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2 && \text{Bias Squared} \\ &\quad + \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) && \text{Variance} \\ &\quad + \text{Var}(y | \vec{x}) && \text{Irreducible Error}\end{aligned}$$

Bias-Variance Decomposition

$$\epsilon(\vec{x}, f, L) = \underbrace{E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}]}_{\text{Expected Error}}$$

Testing error averaged over training datasets compared to noisy versions of the label

$$= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2 \quad \text{Bias Squared}$$

$$+ \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) \quad \text{Variance}$$

$$+ \text{Var}(y | \vec{x}) \quad \text{Irreducible Error}$$

Bias-Variance Decomposition

$$\epsilon(\vec{x}, f, L) = E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}] \quad \text{Expected Error}$$

$$= \left(\underbrace{E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}]}_{\text{Prediction of the trained model on testing example, averaged over training datasets}} - E_y[y | \vec{x}] \right)^2 \quad \text{Bias Squared}$$

$$+ \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) \quad \text{Variance}$$

$$+ \text{Var}(y | \vec{x}) \quad \text{Irreducible Error}$$

Bias-Variance Decomposition

$$\epsilon(\vec{x}, f, L) = E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}] \quad \text{Expected Error}$$

$$= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - \underbrace{E_y[y | \vec{x}]} \right)^2 \quad \text{Bias Squared}$$

Average of noisy versions of the label for the testing example

$$+ \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) \quad \text{Variance}$$

$$+ \text{Var}(y | \vec{x}) \quad \text{Irreducible Error}$$

Bias-Variance Decomposition

$$\epsilon(\vec{x}, f, L) = E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}] \quad \text{Expected Error}$$

$$= \underbrace{\left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2}_{\text{Squared difference between average prediction and average label}} \quad \text{Bias Squared}$$

$$+ \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) \quad \text{Variance}$$

$$+ \text{Var}(y | \vec{x}) \quad \text{Irreducible Error}$$

Bias-Variance Decomposition

$$\epsilon(\vec{x}, f, L) = E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}] \quad \text{Expected Error}$$

$$= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2 \quad \text{Bias Squared}$$

$$+ \underbrace{\text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x})}_{\text{Variance of the prediction on the testing example over different training datasets}} \quad \text{Variance}$$

Variance of the prediction on the testing example over different training datasets

$$+ \text{Var}(y | \vec{x}) \quad \text{Irreducible Error}$$

Bias-Variance Decomposition

$$\epsilon(\vec{x}, f, L) = E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}] \quad \text{Expected Error}$$

$$= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2 \quad \text{Bias Squared}$$

$$+ \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) \quad \text{Variance}$$

$$+ \underbrace{\text{Var}(y | \vec{x})}_{\text{Irreducible Error}}$$

Variance of different noisy versions of the label of the testing example

Bias-Variance Tradeoff

When overfitting:

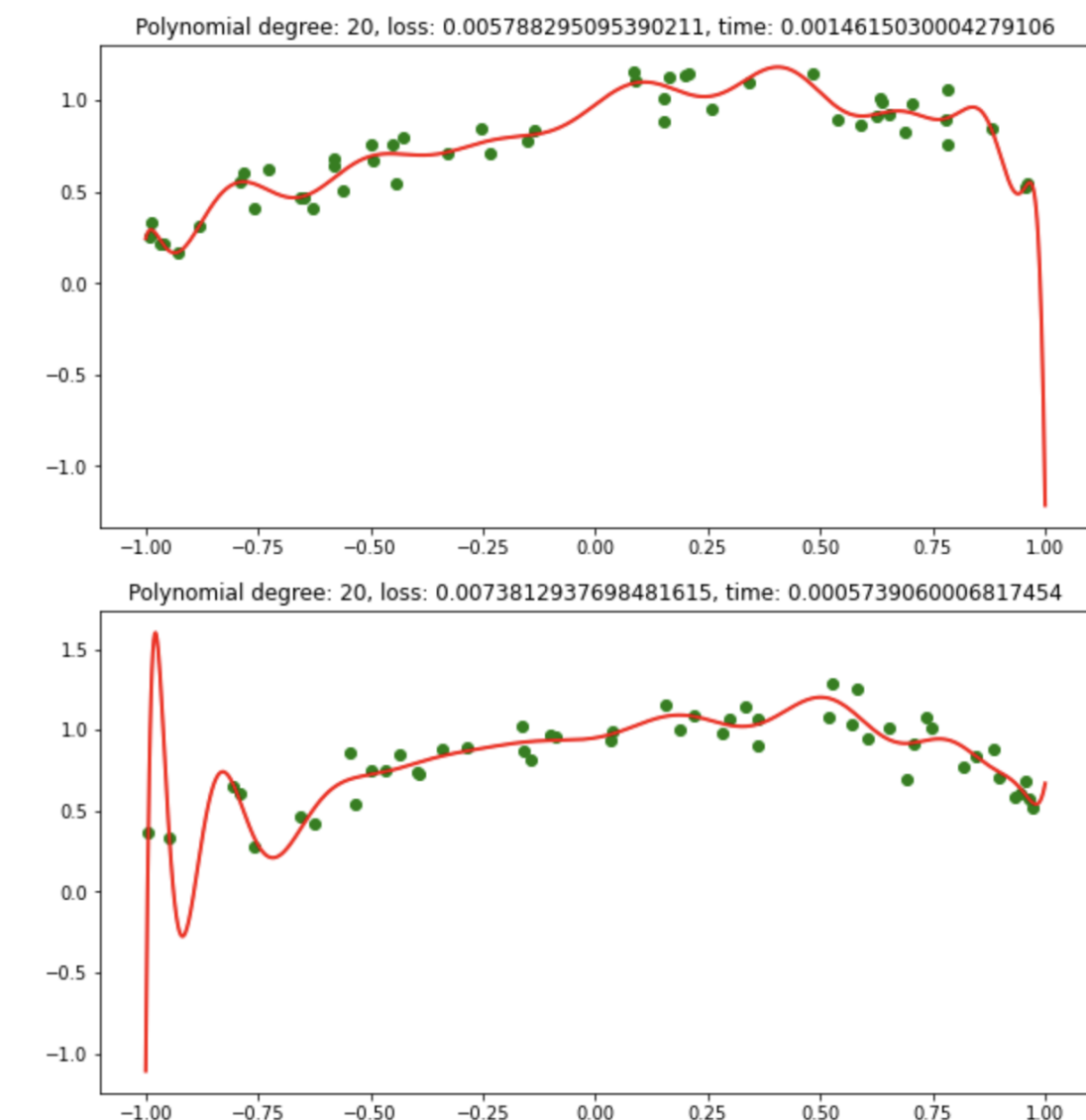
$$\epsilon(\vec{x}, f, L) = E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}] \quad \text{Expected Error}$$

$$= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2$$

$$+ \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) \quad \text{Variance: Large}$$

$$+ \text{Var}(y | \vec{x}) \quad \text{Irreducible Error}$$

Bias Squared: Small



Bias-Variance Tradeoff

When underfitting:

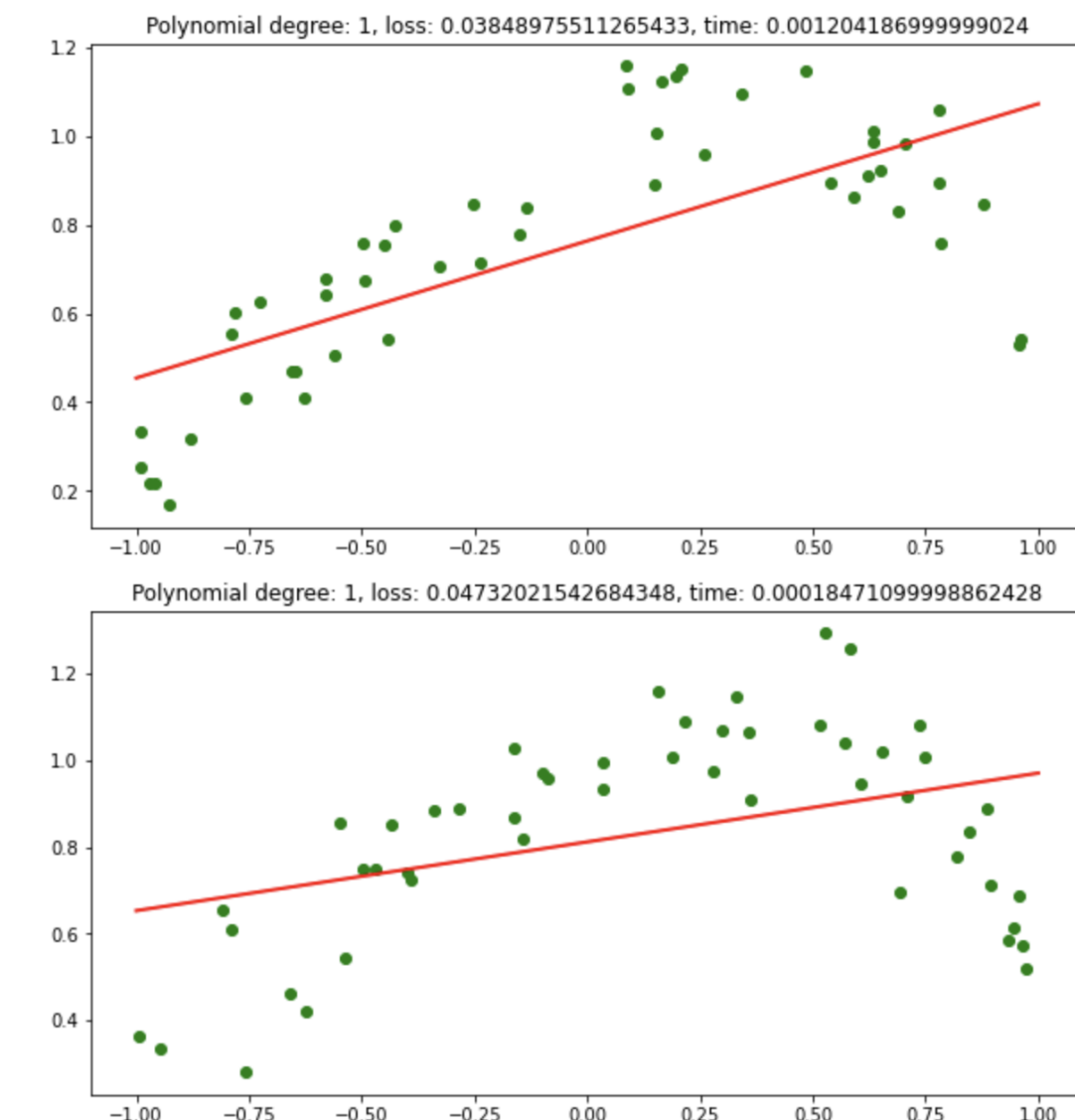
$$\epsilon(\vec{x}, f, L) = E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}] \quad \text{Expected Error}$$

$$= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2$$

$$+ \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) \quad \text{Variance: Small}$$

$$+ \text{Var}(y | \vec{x}) \quad \text{Irreducible Error}$$

Bias Squared: Large



Bias-Variance Tradeoff

When fitted just right:

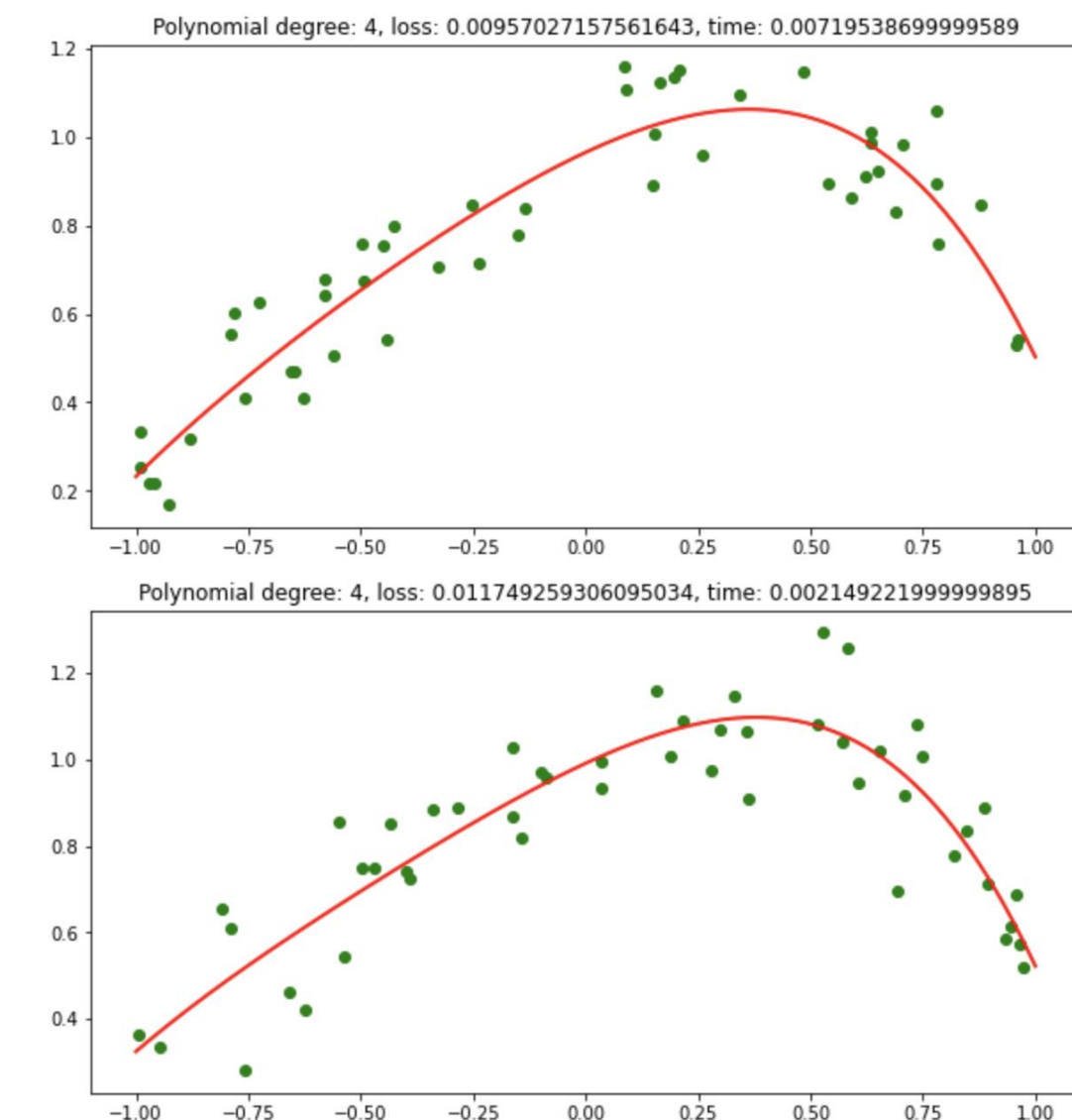
$$\epsilon(\vec{x}, f, L) = E_{y, \mathcal{D}}[(f(\vec{x}; \theta^*(\mathcal{D}, L)) - y)^2 | \vec{x}] \quad \text{Expected Error}$$

$$= \left(E_{\mathcal{D}}[f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}] - E_y[y | \vec{x}] \right)^2$$

$$+ \text{Var}(f(\vec{x}; \theta^*(\mathcal{D}, L)) | \vec{x}) \quad \text{Variance: Small}$$

$$+ \text{Var}(y | \vec{x}) \quad \text{Irreducible Error}$$

Bias Squared: Small



Takeaways

Expected Error on Testing Example = $\text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$

Bias: Difference between average prediction and average label of the testing example

Variance: Variance of the prediction on the testing example over different training datasets

Both bias and variance must be low in order for the expected error to be low

When overfitting, bias is low and variance is high

When underfitting, bias is high and variance is low

Simply achieving low bias (which happens when the model is very expressive) isn't enough!