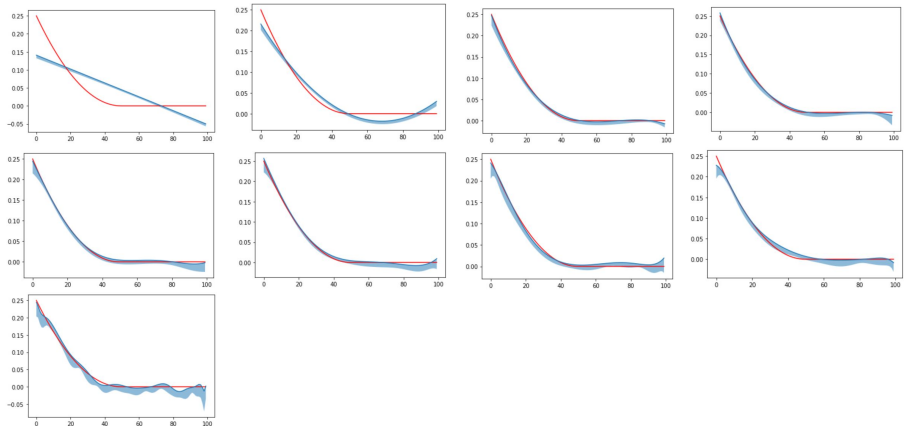


Deep Learning

2.3 Bias-Variance trade-off

Dr. Konda Reddy Mopuri
kmopuri@iittp.ac.in
Dept. of CSE, IIT Tirupati
Aug-Dec 2021

Visualize overfitting



- 1 If we formalize the observations

- ① If we formalize the observations
- ② Let x be fixed, y be the true value associated with it, f^* is what we learned from dataset \mathcal{D} , and $Y = f^*(x)$ is the predicted value

- ① If we formalize the observations
- ② Let x be fixed, y be the true value associated with it, f^* is what we learned from dataset \mathcal{D} , and $Y = f^*(x)$ is the predicted value
- ③ Let's consider that \mathcal{D} is a random variable, then f^* and Y get random

Consider

$$\begin{aligned}
 \mathbb{E}_{\mathcal{D}}((Y - y)^2) &= \mathbb{E}_{\mathcal{D}}(Y^2 - 2Yy + y^2) \\
 &= \mathbb{E}_{\mathcal{D}}(Y^2) - 2\mathbb{E}_{\mathcal{D}}(Y)y + y^2 \\
 &= \mathbb{E}_{\mathcal{D}}(Y^2) - \mathbb{E}_{\mathcal{D}}(Y)^2 + \mathbb{E}_{\mathcal{D}}(Y)^2 + 2\mathbb{E}_{\mathcal{D}}(Y)y + y^2 \\
 &= \mathbb{E}_{\mathcal{D}}(Y^2) - \mathbb{E}_{\mathcal{D}}(Y)^2 + \mathbb{E}_{\mathcal{D}}(Y)^2 + 2\mathbb{E}_{\mathcal{D}}(Y)y + y^2 \\
 &= (\mathbb{E}_{\mathcal{D}}(Y) - y)^2 + \mathbb{V}_{\mathcal{D}}(Y)
 \end{aligned}$$

Bias-Variance Decomposition

$$\begin{aligned}
 \mathbb{E}_{\mathcal{D}}((Y - y)^2) &= \mathbb{E}_{\mathcal{D}}(Y^2 - 2Yy + y^2) \\
 &= \mathbb{E}_{\mathcal{D}}(Y^2) - 2\mathbb{E}_{\mathcal{D}}(Y)y + y^2 \\
 &= \mathbb{E}_{\mathcal{D}}(Y^2) - \mathbb{E}_{\mathcal{D}}(Y)^2 + \mathbb{E}_{\mathcal{D}}(Y)^2 + 2\mathbb{E}_{\mathcal{D}}(Y)y + y^2 \\
 &= \mathbb{E}_{\mathcal{D}}(Y^2) - \mathbb{E}_{\mathcal{D}}(Y)^2 + \mathbb{E}_{\mathcal{D}}(Y)^2 + 2\mathbb{E}_{\mathcal{D}}(Y)y + y^2 \\
 &= (\mathbb{E}_{\mathcal{D}}(Y) - y)^2 + \mathbb{V}_{\mathcal{D}}(Y)
 \end{aligned}$$

① This is known as Bias-Variance decomposition

Bias-Variance Decomposition

$$\begin{aligned}
 \mathbb{E}_{\mathcal{D}}((Y - y)^2) &= \mathbb{E}_{\mathcal{D}}(Y^2 - 2Yy + y^2) \\
 &= \mathbb{E}_{\mathcal{D}}(Y^2) - 2\mathbb{E}_{\mathcal{D}}(Y)y + y^2 \\
 &= \mathbb{E}_{\mathcal{D}}(Y^2) - \mathbb{E}_{\mathcal{D}}(Y)^2 + \mathbb{E}_{\mathcal{D}}(Y)^2 + 2\mathbb{E}_{\mathcal{D}}(Y)y + y^2 \\
 &= \mathbb{E}_{\mathcal{D}}(Y^2) - \mathbb{E}_{\mathcal{D}}(Y)^2 + \mathbb{E}_{\mathcal{D}}(Y)^2 + 2\mathbb{E}_{\mathcal{D}}(Y)y + y^2 \\
 &= (\mathbb{E}_{\mathcal{D}}(Y) - y)^2 + \mathbb{V}_{\mathcal{D}}(Y)
 \end{aligned}$$

- ① This is known as Bias-Variance decomposition
- ② Bias term quantifies how much the model fits the data on average

Bias-Variance Decomposition

$$\begin{aligned}
 \mathbb{E}_{\mathcal{D}}((Y - y)^2) &= \mathbb{E}_{\mathcal{D}}(Y^2 - 2Yy + y^2) \\
 &= \mathbb{E}_{\mathcal{D}}(Y^2) - 2\mathbb{E}_{\mathcal{D}}(Y)y + y^2 \\
 &= \mathbb{E}_{\mathcal{D}}(Y^2) - \mathbb{E}_{\mathcal{D}}(Y)^2 + \mathbb{E}_{\mathcal{D}}(Y)^2 + 2\mathbb{E}_{\mathcal{D}}(Y)y + y^2 \\
 &= \mathbb{E}_{\mathcal{D}}(Y^2) - \mathbb{E}_{\mathcal{D}}(Y)^2 + \mathbb{E}_{\mathcal{D}}(Y)^2 + 2\mathbb{E}_{\mathcal{D}}(Y)y + y^2 \\
 &= (\mathbb{E}_{\mathcal{D}}(Y) - y)^2 + \mathbb{V}_{\mathcal{D}}(Y)
 \end{aligned}$$

- ① This is known as Bias-Variance decomposition
- ② Bias term quantifies how much the model fits the data on average
- ③ Variance term quantifies how much the model changes across datasets

Bias-Variance Trade-off

- ① Reducing the capacity makes f^* fit the data less on average, which increases the bias term

Bias-Variance Trade-off

- ① Reducing the capacity makes f^* fit the data less on average, which increases the bias term
- ② Increasing the capacity makes f^* vary a lot with the training data, which increases the variance term