

## Solution for Assignment 1

### 1. Polynomial Regression

#### 1.1. Code snippet

The completed code has a function `create_dataset` that creates a dataset with sample size = `sample_size` within the range `x_range` for the polynomial whose coefficients are defined by `w_star` with some added noise with normal distribution defined by `sigma = sigma`.

#### 1.2. 2D Scatter plots for training and validation set

Upon using the above function with the following parameters a training and validation set was generated.

#### 1.3. bias parameter in `torch.nn.linear`

Bias in Linear Regression is just like an intercept added in a linear equation. It is an additional parameter in the Neural Network which is used to adjust the output along with the weighted sum of the inputs to the neuron. Moreover, bias value allows you to shift the activation function to either right or left. An example of changing bias can be seen in figure: 1

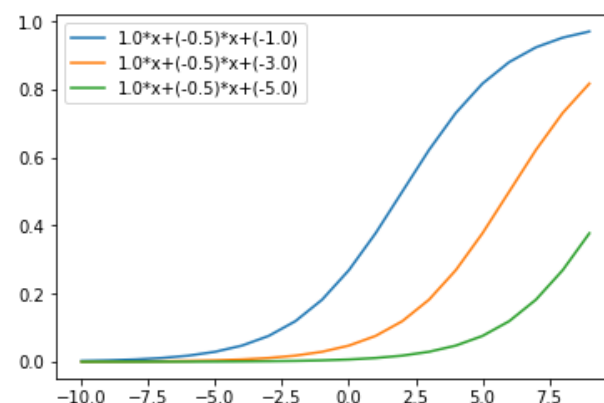


Figure 1: Bias

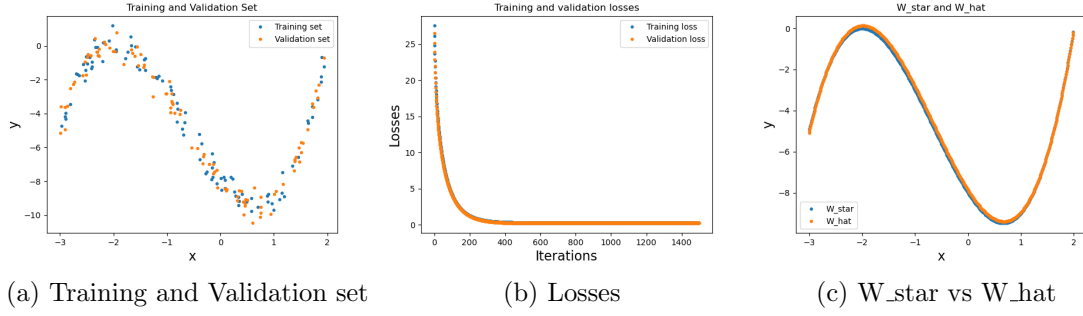


Figure 2: Linear Regression

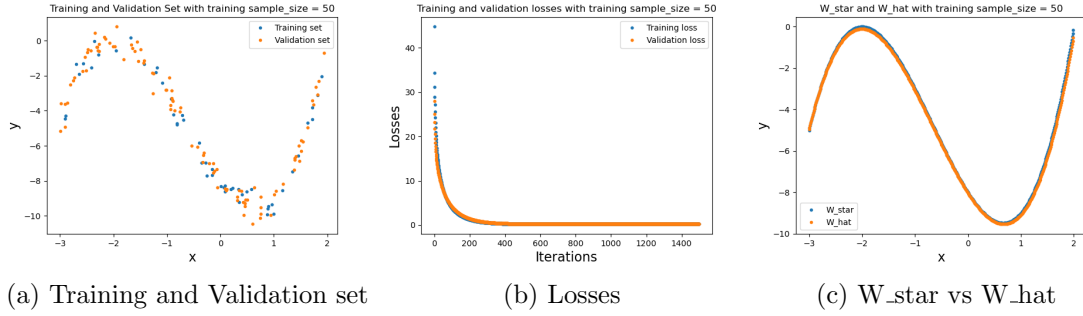


Figure 3: Training sample\_size = 50

## 1.4. Linear Regression

Using `nn.linear`, `optim.SGD` methods from the `pyTorch` library to build a single layer neural network to find the parameters in  $w\_star$  using Stochastic Gradient Descent optimization. The suitable `learning_rate` and `max_iterations` was determined to be 0.01 and 1500 respectively.

In the plot 2b, we can see that the training and validation loss steadily drops as more iterations go by. This is expected of an SGD optimizer where we are minimising loss by finding the negative gradient of the loss function. It can be noted that the max number of iterations can actually be reduced to the range of 600-800 since the optimization problem seemed to have converged.

In plot 2c, the plot for the polynomial defined by input vector  $w\_star$  and result of linear regression  $w\_hat$  with constant  $w_0$  can be seen. Visually speaking, the results are pretty accurate and mathematically the losses have been minimized to a magnitude of 0.25.

## 1.5. Reducing training data

Reducing the training data to 50, 10 and 5 leads to interesting results. As the number of training samples keep dropping the training data of course becomes more sparse on the plot with training and validation set.

### 1.5.1. Training Data = 50 samples

In decreasing the training samples to 50, the loss vs iteration graph in 3b doesn't look very different from the original case. But that is visually of course, mathematically we're seeing a slightly underfit model because there's not enough training data to learn from. The validation losses are also reducing so its safe to say that the general trend has been learnt pretty well, evident from the  $W\_star$  vs  $W\_hat$  where the equation line and prediction line are overlapping.

### 1.5.2. Training Data = 10 samples

In decreasing the training samples to 10, the loss vs iteration graph in 4b looks vastly different. The clear signs of overfit model is visible because the validation losses shoot up and the start reducing.

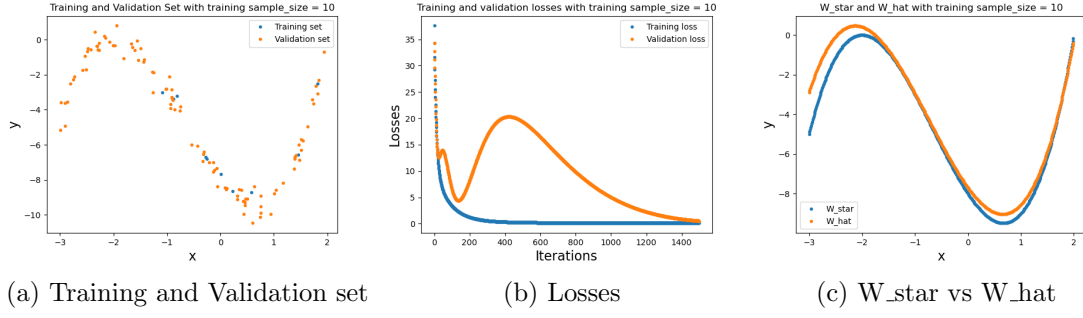


Figure 4: Training sample\_size = 10

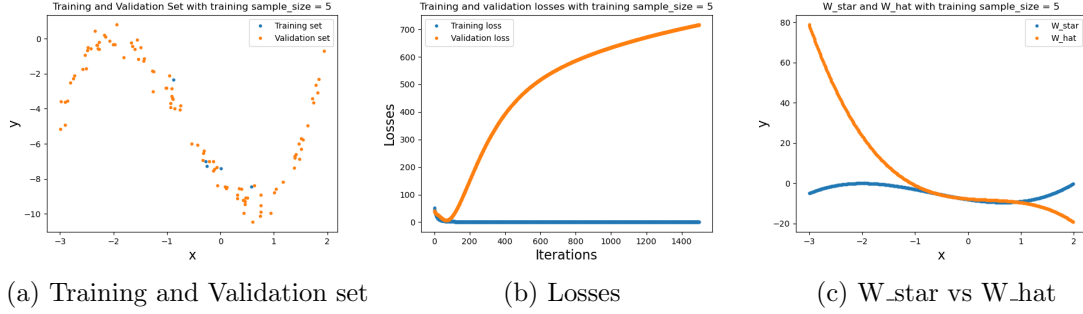


Figure 5: Training sample\_size = 5

Perhaps upon increasing the number of iterations, we could eventually arrive at a solution where the validation and training losses are as low as with more number of training samples. The overfit behavior is more evident in the  $W_{star}$  vs  $W_{hat}$  where the prediction line is not in perfect overlap with the equation line.

### 1.5.3. Training Data = 5 samples

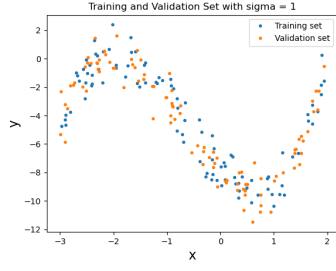
In decreasing the training samples to 5, the loss vs iteration graph in 5b depicts awful performance from the model. The signs of an incorrect model are clearly visible with exponentially increasing validation losses while the training loss is ever reducing. This is evidence that because of low number of training samples, the model is overfit to those 5 data points. Hence the training losses are practically 0 but validation losses are extremely high. The incorrect model is more evident in the  $W_{star}$  vs  $W_{hat}$  where the prediction line is nowhere close to what the equation line is. It is evident that the model overfit to a completely different polynomial because the low training data and high iterations reinforced a bad generalisation in the model.

## 1.6. Increasing $\sigma$

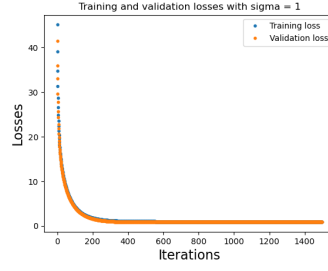
Increasing the sigma when generating the training and validation set introduces a lot more noise when generating the datasets. This would lead to an underfit model since the variation in data is too high for the neural network to learn. So the predicted  $w_{hat}$  is vastly different from the actual  $w_{star}$  and the variance keeps increasing as we increase sigma. Another peculiar behaviour observed was that the training and validation loss have odd behaviors as we increase sigma where the validation error is lower than the training error. This could be because of the order of noise added is the same as the order of values of  $y$ .

## 1.7. Bonus Task

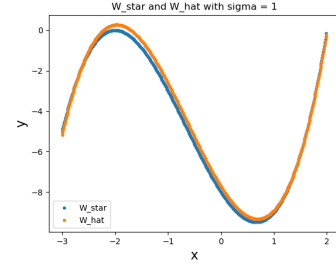
Upon increasing the degree of the polynomial from 3 to 4 with 10 data points for training, the combined effect of higher complexity and lower training data points is evident from the behavior



(a) Training and Validation set

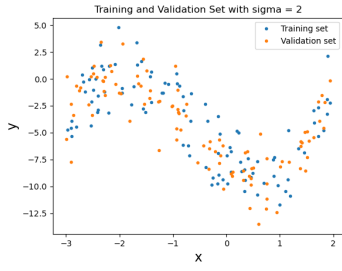


(b) Losses

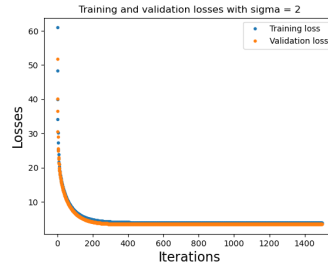


(c)  $W_{\text{star}}$  vs  $W_{\text{hat}}$

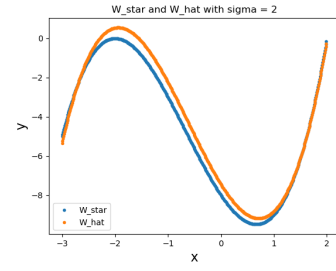
Figure 6:  $\sigma = 1$



(a) Training and Validation set



(b) Losses

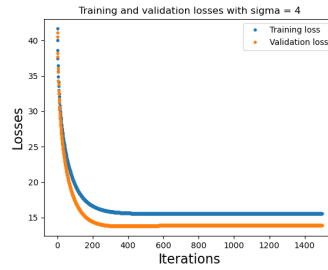


(c)  $W_{\text{star}}$  vs  $W_{\text{hat}}$

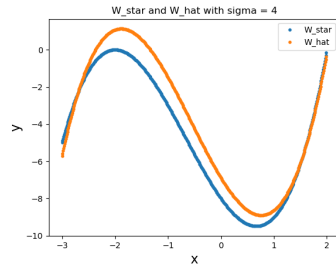
Figure 7:  $\sigma = 2$



(a) Training and Validation set



(b) Losses

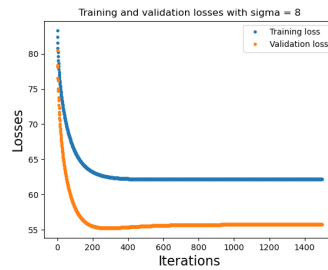


(c)  $W_{\text{star}}$  vs  $W_{\text{hat}}$

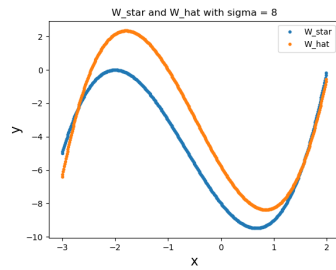
Figure 8:  $\sigma = 4$



(a) Training and Validation set



(b) Losses



(c)  $W_{\text{star}}$  vs  $W_{\text{hat}}$

Figure 9:  $\sigma = 8$

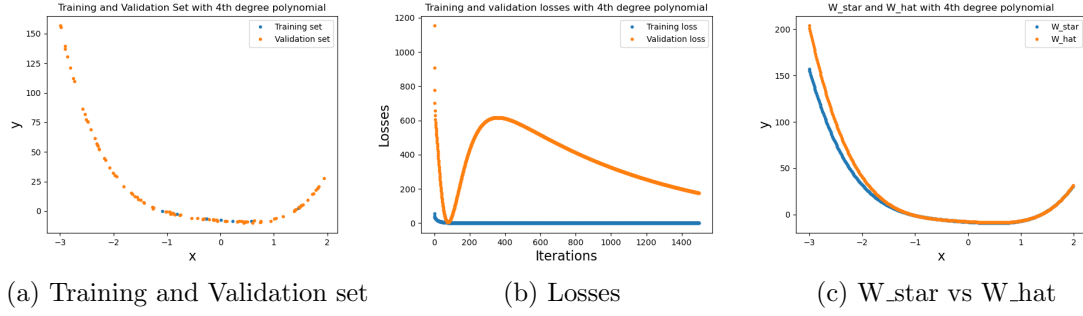


Figure 10: 4th Degree polynomial

observed in the graphs on figure: 10. The magnitude of the loss metrics for predicting a 4th degree polynomial with 10 training data points is an order higher than the case with 3rd. degree polynomial. Same high order characteristics can be observed in the the difference between graph of  $w_{star}$  and  $w_{hat}$ .

The original 3rd order polynomial:

$$p(x) = x^3 + 2x^2 - 4x - 8 \quad (1)$$

The 4th order polynomial:

$$p(x) = 2x^4 + x^3 + 2x^2 - 4x - 8 \quad (2)$$

## 2. Questions

### 2.1. Local and Global Minima

Local minima is the point in a range for which the value of the polynomial is lower than that of the neighboring points. If in addition to this the value of polynomial is infimum of the whole range, then the points is considered a global minima.

### 2.2. Underfitting and overfitting

When we say a model is underfit, that means it is unable to learn the patterns in the data properly. An underfit model doesn't fully learn each and every example in the dataset. When we say a model is overfit, that means it is over trained on the data such that, it even learns the noise from it. An overfit model learns each and every example so perfectly that it misclassifies an unseen/new example.

### 2.3. Fixing overfit and underfit models

To fix an underfit model many remedies can be used, one can be to increase the complexity of the model for a better fit. For example, using a 3rd order polynomial as the regression line instead of 2nd order polynomial. Vice versa, to fix an overfit model a simpler model should be used for a better fit.