

CS578 Statistical Machine Learning  
02/08/2021 Lecture Notes: Linear Regression

Site Bai

## 1 The Linear Regression Model

Given an input  $[\mathbf{x}_1^\top, \mathbf{x}_2^\top, \dots, \mathbf{x}_N^\top]^\top \in \mathbb{R}^{N \times k}$ ,  $N$  is the number of input data points and  $k$  is the number of features we observe for each data point. We want to find a mapping between the input and the target variables  $\mathbf{y} = [y_1, y_2, \dots, y_N]^\top \in \mathbb{R}^N$ , as shown in Figure 1. For each input data and its corresponding target variable, the linear regression model assume the target variable can be represented as a linear combination of the input with some Gaussian noise, i.e.

$$y_i = w_0 + \sum_{j=1}^k w_j x_{ij} + \epsilon, \quad (1)$$

in which  $i = 1, 2, \dots, N$ ,  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ . To rewrite (1) in a matrix format,

$$y_i = \mathbf{w}^\top \begin{bmatrix} 1 \\ \mathbf{x}_i \end{bmatrix} + \epsilon. \quad (2)$$

$\mathbf{w} \in \mathbb{R}^{k+1}$ ,  $\mathbf{x}_i \in \mathbb{R}^k$ . Denote  $\mathbf{x}'_i = [1, \mathbf{x}_i]^\top$ , then  $y_i = \mathbf{w}^\top \mathbf{x}'_i + \epsilon$ , and  $\mathbb{P}(y_i | \mathbf{x}'_i, \mathbf{w}) \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}'_i, \sigma^2)$ .

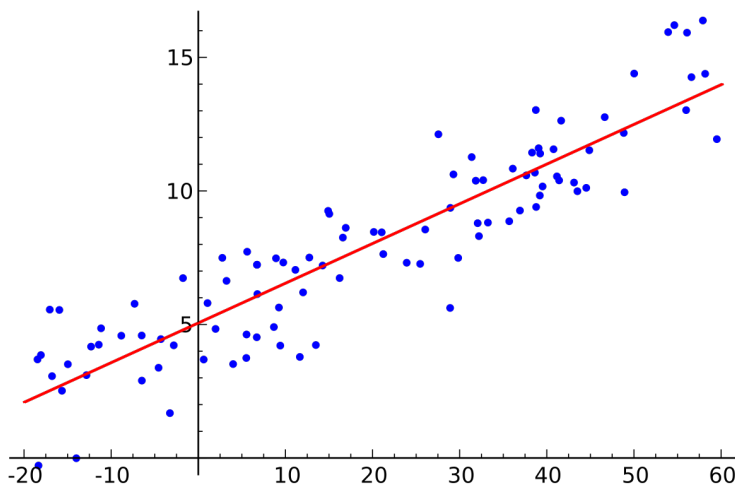


Figure 1: Linear Regression, figure from Wikipedia.

## 2 Maximum Likelihood Estimation

To find a proper mapping from input to target, we need to obtain the parameter  $\mathbf{w}$  that fits the data, and this can be done by Maximum Likelihood Estimation (MLE). We have

$$\mathbb{P}(y_i|\mathbf{x}'_i, \mathbf{w}) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(y_i - \mathbf{w}^\top \mathbf{x}'_i)^2}{2\sigma^2} \right\} \quad (3)$$

Thus, the log likelihood function can be written as

$$\mathcal{L}(\mathbf{w}) = \log \prod_{i=1}^N \mathbb{P}(y_i|\mathbf{x}'_i, \mathbf{w}) \quad (4)$$

$$= \sum_{i=1}^N \log \mathbb{P}(y_i|\mathbf{x}'_i, \mathbf{w}) \quad (5)$$

$$= \sum_{i=1}^N \left( \log \frac{1}{\sqrt{2\pi}\sigma} + \log \exp \left\{ -\frac{(y_i - \mathbf{w}^\top \mathbf{x}'_i)^2}{2\sigma^2} \right\} \right) \quad (6)$$

$$= \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}'_i)^2 \quad (7)$$

And to maximize the likelihood,

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} \mathcal{L}(\mathbf{w}) \quad (8)$$

$$= \arg \min_{\mathbf{w}} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}'_i)^2 \quad (9)$$

$$= \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad (10)$$

$$= \arg \min_{\mathbf{w}} \text{RSS}(\mathbf{w}), \quad (11)$$

in which  $\mathbf{X} = [\mathbf{x}'_1{}^\top, \mathbf{x}'_2{}^\top, \dots, \mathbf{x}'_N{}^\top]^\top \in \mathbb{R}^{N \times (k+1)}$ , RSS denotes the residual sum-of-squares.

## 3 Two Ways of Computing $\hat{\mathbf{w}}$

### 3.1 Derivative of $\text{RSS}(\mathbf{w})$

We can compute the derivative of  $\text{RSS}(\mathbf{w})$  w.r.t.  $\mathbf{w}$  to obtain the optimal  $\mathbf{w}$ .

$$\frac{\partial \text{RSS}(\mathbf{w})}{\partial \mathbf{w}} = -2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}). \quad (12)$$

Then set the derivative to zero,

$$-2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) = 0. \quad (13)$$

We now obtain the solution

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (14)$$

### 3.2 Reformatting the Quadratic Form

$$\text{RSS}(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad (15)$$

$$= (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) \quad (16)$$

$$= \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{y}^\top \mathbf{X} \mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y} \quad (17)$$

$$= \mathbf{w}^\top (\mathbf{X}^\top \mathbf{X}) \mathbf{w} - \mathbf{y}^\top \mathbf{X} \mathbf{w} - \underbrace{\mathbf{w}^\top (\mathbf{X}^\top \mathbf{X}) (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}_{\mathbf{I}} \quad (18)$$

$$+ \underbrace{\mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}}_0$$

$$= (\mathbf{w}^\top (\mathbf{X}^\top \mathbf{X}) - \mathbf{y}^\top \mathbf{X}) \mathbf{w} - (\mathbf{w}^\top (\mathbf{X}^\top \mathbf{X}) - \mathbf{y}^\top \mathbf{X}) (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (19)$$

$$- \underbrace{\mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}}_{\text{Constant}}$$

$$= (\mathbf{w}^\top (\mathbf{X}^\top \mathbf{X}) - \mathbf{y}^\top \mathbf{X}) \left( \mathbf{w} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \right) + \text{Constant} \quad (20)$$

$$= \left( \mathbf{w}^\top (\mathbf{X}^\top \mathbf{X}) - \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{X}) \right) \left( \mathbf{w} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \right) + \text{Constant} \quad (21)$$

$$= \left( \mathbf{w}^\top - \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \right) (\mathbf{X}^\top \mathbf{X}) \left( \mathbf{w} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \right) + \text{Constant} \quad (22)$$

$$= \left( \underbrace{\mathbf{w} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}_{\mathbf{U}} \right)^\top (\mathbf{X}^\top \mathbf{X}) \left( \underbrace{\mathbf{w} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}_{\mathbf{U}} \right) + \text{Constant} \quad (23)$$

$$= \mathbf{U}^\top \mathbf{X}^\top \mathbf{X} \mathbf{U} + \text{Constant} \quad (24)$$

$$= (\mathbf{X} \mathbf{U})^\top (\mathbf{X} \mathbf{U}) + \text{Constant} \quad (25)$$

$$= \|\mathbf{X} \mathbf{U}\|_2^2 + \text{Constant} \quad (26)$$

$\text{RSS}(\mathbf{w})$  is minimized when  $\|\mathbf{X} \mathbf{U}\|_2^2 = 0$ , i.e.  $\mathbf{U} = 0$ . Thus, we have

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (27)$$

## 4 About $(\mathbf{X}^\top \mathbf{X})^{-1}$

This section poses the question that what if  $\mathbf{X}^\top \mathbf{X}$  is not invertible. We know that  $\mathbf{X} \in \mathbb{R}^{N \times (k+1)}$ , thus  $\mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{(k+1) \times (k+1)}$  is a real symmetric  $(k+1) \times (k+1)$  matrix. According to the symmetric eigenvalue decomposition of  $\mathbf{X}^\top \mathbf{X}$ ,

$$\mathbf{X}^\top \mathbf{X} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top, \quad (28)$$

in which  $\mathbf{Q} \in \mathbb{R}^{(k+1) \times (k+1)}$  is orthogonal, and  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_{k+1})$ .

If  $\mathbf{X}^\top \mathbf{X}$  is invertible,

$$(\mathbf{X}^\top \mathbf{X})^{-1} = \mathbf{Q} \mathbf{\Lambda}^{-1} \mathbf{Q}^\top, \quad (29)$$

in which  $\mathbf{\Lambda}^{-1} = \mathbf{diag}\left(\frac{1}{\lambda_1}, \frac{1}{\lambda_2}, \dots, \frac{1}{\lambda_{k+1}}\right)$ . If  $\mathbf{X}^\top \mathbf{X}$  is not invertible, i.e.  $\exists \lambda_i \in \{\lambda_1, \lambda_2, \dots, \lambda_{k+1}\}$  that  $\lambda_i = 0$ , in practice, we add some small  $\lambda > 0$  so that  $\mathbf{\Lambda}^{-1} = \mathbf{diag}\left(\frac{1}{\lambda_1 + \lambda}, \frac{1}{\lambda_2 + \lambda}, \dots, \frac{1}{\lambda_{k+1} + \lambda}\right)$ , i.e.

$$\hat{\mathbf{w}} = \left(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}\right)^{-1} \mathbf{X}^\top \mathbf{y}. \quad (30)$$

## 5 Parametric v.s. Non-Parametric Models

As shown in previous sections, we can notice that the model of Linear Regression is based on computing a parameter  $\mathbf{w}$ , i.e. the mapping from the input to the target variables is determined by this parameter  $\mathbf{w}$ . We call this kind of Machine Learning models the Parametric Models. In comparison, those without the need to embed the mapping into parameters are called Non-Parametric Models. K-Nearest Neighbor is one example of Non-Parametric Models.

## 6 Regression with Non-linear Basis

The assumption that the target variables can be represented by the linear combination of the input data features with Gaussian noise may not always hold. For some obvious non-linear samples as shown in Figure 2, we can fit the data by regression with non-linear basis. In the setting of regression with non-linear basis,

$$y_i = \mathbf{w}^\top \phi(\mathbf{x}_i). \quad (31)$$

$\phi(\mathbf{x}_i)$  is some non-linear function of  $\mathbf{x}_i$ , polynomials to name an example. Then,

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N \left(y_i - \mathbf{w}^\top \phi(\mathbf{x}_i)\right)^2 \quad (32)$$

$$= \|\mathbf{y} - \mathbf{\Phi} \mathbf{w}\|_2^2, \quad (33)$$

in which  $\mathbf{\Phi} = (\phi^\top(\mathbf{x}_1), \phi^\top(\mathbf{x}_2), \dots, \phi^\top(\mathbf{x}_N))^\top$ . Following similar derivation in Section 3, we have

$$\hat{\mathbf{w}} = \left(\mathbf{\Phi}^\top \mathbf{\Phi}\right)^{-1} \mathbf{\Phi}^\top \mathbf{y}. \quad (34)$$

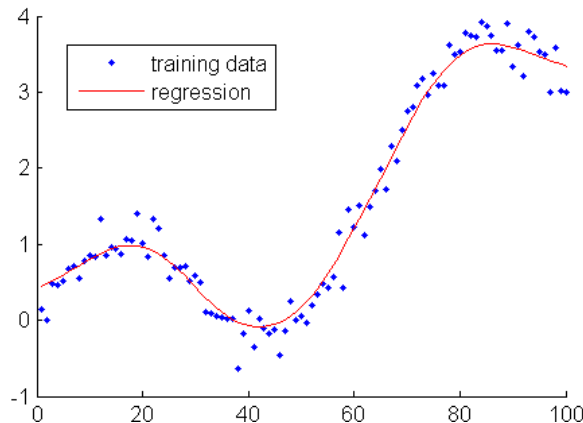


Figure 2: Non-Linear Regression, figure from Wikipedia.

## 7 Underfitting v.s. Overfitting

The non-linear basis of regression can be of different levels of degree or complexity. The bad choice of the complexity may lead to the problem of underfitting or overfitting. For polynomial basis, the highest order of the polynomial is the obvious representation of model complexity. As shown in Figure 3, different degrees lead to different fitting results.

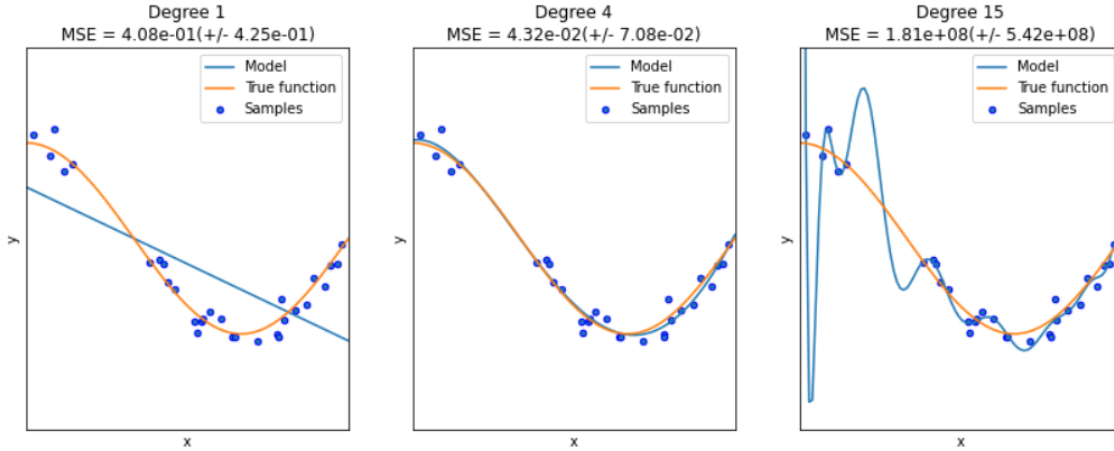


Figure 3: Underfitting v.s. Overfitting, figure from Data Science Foundation.

If the degree is too low, then the basis given is unable to fit the given sample; that is to say the bias of underfitting is high. And the result would not change much no matter how many attempts are made because it reaches the limit of the representation capability of low degree functions. Thus we can say the variance of overfitting is low.

Overfitting, on the other hand, would lure the function to pass every given data points, resulting in a possibly very low bias. However, if the given data changes slightly, possibly due to some random noise, the result may change dramatically, which leads to a high variance.

To find out the best complexity of the model, in practice we usually observe the error on both the training set and the evaluation set as the complexity increases. The common cases follow the trend shown in Figure 4.

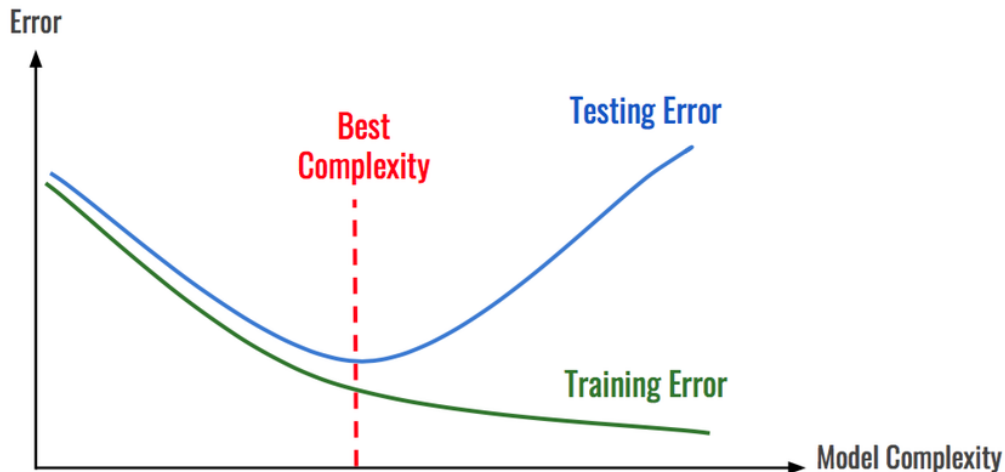


Figure 4: Best Choice of Model Complexity, figure by Robert Martin-Short.

## 8 Bias-Variance Trade-off

This section studies the bias-variance trade-off in regression models, which can be also regarded as a theoretical and quantitative study of the under/over-fitting problem.

We show that the error of regression is composed by a variance term and a bias term (possibly another noise). We assume that  $\exists h(\mathbf{x}_i)$  that is the optimal prediction function for predicting the target variable,

$$h(\mathbf{x}) = \mathbb{E}_y [\mathbb{P}(y | \mathbf{x})] \quad (35)$$

and

$$y_i = h(\mathbf{x}_i) + \epsilon, \quad (36)$$

$\epsilon \sim \mathcal{N}(0, \sigma^2)$ . It's important to notice that for any input  $\mathbf{x}_i$ ,  $h(\mathbf{x}_i)$  is deterministic. Our prediction based on some dataset  $\mathcal{D}$  given the linear model in (1) is denoted by  $f_{\mathcal{D}}(\mathbf{x}_i)$ . We can write the expectation of the loss function (Expected Square Loss) as follows:

$$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}(\mathbf{x}, y)} [L] = \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}(\mathbf{x}, y)} \left[ (y - f_{\mathcal{D}}(\mathbf{x}))^2 \right]. \quad (37)$$

The expectation of the error over different datasets is

$$\mathbb{E}_{\mathcal{D}, (\mathbf{x}, y)} [L] = \mathbb{E}_{\mathcal{D}} \mathbb{E}_{(\mathbf{x}, y) \sim p(\mathbf{x}, y)} \left[ (y - f_{\mathcal{D}}(\mathbf{x}))^2 \right] \quad (38)$$

$$= \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}} \left[ (h(\mathbf{x}) + \epsilon - f_{\mathcal{D}}(\mathbf{x}))^2 \right] \quad (39)$$

$$= \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}} \left[ \left( \underbrace{h(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})]}_{\mathbf{A}} + \underbrace{\mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] - f_{\mathcal{D}}(\mathbf{x})}_{\mathbf{B}} + \epsilon \right)^2 \right] \quad (40)$$

$$= \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}} \left[ \mathbf{A}^2 + \mathbf{B}^2 + \epsilon^2 + 2\epsilon\mathbf{A} + 2\epsilon\mathbf{B} + 2\mathbf{A}\mathbf{B} \right] \quad (41)$$

$$= \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}} \left[ \left( h(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] \right)^2 \right] + \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}} \left[ \left( \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] - f_{\mathcal{D}}(\mathbf{x}) \right)^2 \right] + \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}} [\epsilon^2] \quad (42)$$

$$+ 2 \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}} \left[ \epsilon \left( h(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] \right) \right] + 2 \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}} \left[ \epsilon \left( \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] - f_{\mathcal{D}}(\mathbf{x}) \right) \right]$$

$$+ 2 \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}} \left[ \left( h(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] \right) \left( \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] - f_{\mathcal{D}}(\mathbf{x}) \right) \right]$$

$$= \left( h(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] \right)^2 + \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}} \left[ \left( \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] - f_{\mathcal{D}}(\mathbf{x}) \right)^2 \right] + \mathbb{E}[\epsilon^2] \quad (43)$$

$$+ 2 \left( h(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] \right) \mathbb{E}[\epsilon] + 2 \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}} \left[ \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] - f_{\mathcal{D}}(\mathbf{x}) \right] \mathbb{E}[\epsilon]$$

$$+ 2 \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}} \left[ \left( h(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] \right) \left( \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] - f_{\mathcal{D}}(\mathbf{x}) \right) \right].$$

Given that  $\mathbb{E}[\epsilon] = 0$ ,

$$\mathbb{E}_{\mathcal{D},(\mathbf{x},y)}[L] = \left( h(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] \right)^2 + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathcal{D}} \left[ \left( \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] - f_{\mathcal{D}}(\mathbf{x}) \right)^2 \right] + \text{Var}[\epsilon] \quad (44)$$

$$+ 2 \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathcal{D}} \left[ \underbrace{\left( h(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] \right)}_{\text{Constant}} \left( \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] - f_{\mathcal{D}}(\mathbf{x}) \right) \right] \\ = \left( h(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] \right)^2 + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathcal{D}} \left[ \left( \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] - f_{\mathcal{D}}(\mathbf{x}) \right)^2 \right] + \text{Var}[\epsilon] \quad (45)$$

$$+ 2 \left( h(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] \right) \underbrace{\mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathcal{D}} \left[ \left( \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] - f_{\mathcal{D}}(\mathbf{x}) \right) \right]}_0 \\ = \underbrace{\left( h(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] \right)^2}_{\text{Bias}} + \underbrace{\mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathcal{D}} \left[ \left( f_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})] \right)^2 \right]}_{\text{Variance}} + \underbrace{\sigma^2}_{\text{Noise}}. \quad (46)$$

Thus, the Mean Squared Error can be decomposed into a bias term, a variance term, and a noise term, forming a trade-off relationship, i.e. there's no way one can simultaneously minimize the two sources of error.

## 9 Measures to Prevent Overfitting: Regression with Regularization

Measures can be taken to prevent overfitting:

- Train the model on larger datasets;
- Control model complexity;
- Add a regularizer to the regression model.

We can add a regularization term to the regression model to control the model complexity and get

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \|\mathbf{y} - \Phi \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2. \quad (47)$$

To solve this model following similar ways in Section 3, we get

$$\hat{\mathbf{w}} = \left( \Phi^\top \Phi + \lambda \mathbf{I} \right)^{-1} \Phi^\top \mathbf{y}. \quad (48)$$