

CIS 419/519: Homework 2

Jiatong Sun

Although the solutions are entirely my own, I consulted with the following people and sources while working on this homework: *JunfanPan*, *ZhuoyuHe*, *YuchenSun*, *ChangLiu*, *YihangXu*, *YupengLi*
<https://machinelearningmastery.com/understand-the-dynamics-of-learning-rate-on-deep-learning-neural-networks/>
https://en.wikipedia.org/wiki/Learning_rate
<https://machinelearningmastery.com/how-to-tune-algorithm-parameters-with-scikit-learn/>.

1 Gradient Descent

- The implication of the learning rate α_k is to control how big a step should be taken in the gradient descent direction towards the minimum, where a too small α_k may result in a long training time and a too large α_k may lead to an overshooting training process.
- The implications of setting α_k as a function of k is to select an adaptive learning rate based on the training process, since the best step to take can vary as the the training goes gradually towards the minimum and a preset constant α_k may not work well in the whole process.

2 Linear Regression [CIS 519 ONLY]

Since

$$y_i = f(\mathbf{x}_i) + \epsilon_i \quad (1)$$

and

$$\epsilon_i \sim G(0, \sigma^2) \quad (2)$$

We can know that function f is the linear regression function without error

$$f(\mathbf{x}_i) = \theta_0 + \sum_{j=1}^d \theta_j x_{ij} = \sum_{j=0}^d \theta_j x_{ij}, (x_{i0} = 1) \quad (3)$$

or in matrix form

$$f(\mathbf{x}_i) = \mathbf{x}_i \boldsymbol{\theta} \quad (4)$$

where

$$\mathbf{x}_i = [1 \quad x_{i1} \quad \dots \quad x_{ij} \quad \dots \quad x_{id}] \quad (5)$$

$$\boldsymbol{\theta} = [\theta_0 \quad \theta_1 \quad \dots \quad \theta_j \quad \dots \quad \theta_d]^T \quad (6)$$

From the closed form solution, we can write $\boldsymbol{\theta}$ in the following format

$$\boldsymbol{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (7)$$

where \mathbf{X} represents the whole training set and \mathbf{y} represents its label

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_i \\ \vdots \\ \mathbf{x}_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1j} & \dots & x_{1d} \\ 1 & x_{21} & \dots & x_{2j} & \dots & x_{2d} \\ \vdots & \vdots & & \vdots & & \vdots \\ 1 & x_{i1} & \dots & x_{ij} & \dots & x_{id} \\ \vdots & \vdots & & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{nj} & \dots & x_{nd} \end{bmatrix}, \quad (8)$$

$$\mathbf{y} = [y_0 \quad y_1 \quad \dots \quad y_i \quad \dots \quad y_n]^T \quad (9)$$

Let \mathbf{x} be a column vector, which represents *the test data set instead of the training data set*.

From (4), we can write $f(\mathbf{x})$ in the following format

$$f(\mathbf{x}) = h_{\theta}(\mathbf{x}) = \mathbf{x}\theta \quad (10)$$

where

$$\mathbf{x} = \begin{bmatrix} \mathbf{t}_1 \\ \mathbf{t}_2 \\ \vdots \\ \mathbf{t}_k \\ \vdots \\ \mathbf{t}_m \end{bmatrix} = \begin{bmatrix} 1 & t_{11} & \dots & t_{1j} & \dots & t_{1d} \\ 1 & t_{21} & \dots & t_{2j} & \dots & t_{2d} \\ \vdots & \vdots & & \vdots & & \vdots \\ 1 & t_{k1} & \dots & t_{kj} & \dots & t_{kd} \\ \vdots & \vdots & & \vdots & & \vdots \\ 1 & t_{m1} & \dots & t_{mj} & \dots & t_{md} \end{bmatrix}, \quad (11)$$

Here, we use \mathbf{t}_k and \mathbf{t}_m to replace \mathbf{x}_i and \mathbf{x}_n so we can distinguish the training data and the test data.

From (4) and (7), we get

$$f(\mathbf{x}) = \mathbf{x}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (12)$$

where the dimensions are $\mathbf{x} = \mathbf{x}_{m \times (d+1)}$, $\mathbf{X} = \mathbf{X}_{n \times (d+1)}$, $\mathbf{y} = \mathbf{y}_{n \times 1}$

So $\mathbf{x}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ has the dimension of $[m \times n]$, or

$$\mathbf{L}_{m \times n} \triangleq \mathbf{x}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \begin{bmatrix} L_{10} & L_{11} & \dots & L_{1i} & \dots & L_{1n} \\ L_{20} & L_{21} & \dots & L_{2i} & \dots & L_{2n} \\ \vdots & \vdots & & \vdots & & \vdots \\ L_{k0} & L_{k1} & \dots & L_{ki} & \dots & L_{kn} \\ \vdots & \vdots & & \vdots & & \vdots \\ L_{m0} & L_{m1} & \dots & L_{mi} & \dots & L_{mn} \end{bmatrix} \quad (13)$$

So (12) becomes

$$\begin{aligned}
f(\mathbf{x}) = \mathbf{L}_{m \times n} \mathbf{y}_{n \times 1} &= \begin{bmatrix} L_{10} & L_{11} & \dots & L_{1i} & \dots & L_{1n} \\ L_{20} & L_{21} & \dots & L_{2i} & \dots & L_{2n} \\ \vdots & \vdots & & \vdots & & \vdots \\ L_{k0} & L_{k1} & \dots & L_{ki} & \dots & L_{kn} \\ \vdots & \vdots & & \vdots & & \vdots \\ L_{m0} & L_{m1} & \dots & L_{mi} & \dots & L_{mn} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{bmatrix} \\
&= \begin{bmatrix} L_{10} \\ L_{20} \\ \vdots \\ L_{k0} \\ \vdots \\ L_{m0} \end{bmatrix} y_1 + \begin{bmatrix} L_{11} \\ L_{21} \\ \vdots \\ L_{k1} \\ \vdots \\ L_{m1} \end{bmatrix} y_2 + \dots + \begin{bmatrix} L_{1i} \\ L_{2i} \\ \vdots \\ L_{ki} \\ \vdots \\ y_{mi} \end{bmatrix} y_i + \dots + \begin{bmatrix} L_{1n} \\ L_{2n} \\ \vdots \\ L_{kn} \\ \vdots \\ L_{mn} \end{bmatrix} y_n \\
&= \sum_{i=1}^n \begin{bmatrix} L_{1i} \\ L_{2i} \\ \vdots \\ L_{ki} \\ \vdots \\ L_{mi} \end{bmatrix} y_i = \sum_{i=1}^n l_i(x; X) y_i
\end{aligned} \tag{14}$$

so the conclusion is

$$l_i(x; X) = \begin{bmatrix} L_{1i} \\ L_{2i} \\ \vdots \\ L_{ki} \\ \vdots \\ L_{mi} \end{bmatrix} = \text{the } i^{\text{th}} \text{ column of } \mathbf{x}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \tag{15}$$

Note that all calculations above are based on general conditions.

According to this question, since $x_i \in \mathbb{R}$, we know that $d = 1$ and

$$\mathbf{x}_i = \begin{bmatrix} 1 & x_i \end{bmatrix} \tag{16}$$

$$\boldsymbol{\theta} = [\theta_0 \quad \theta_1]^T \tag{17}$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_i \\ \vdots \\ \mathbf{x}_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_i \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} \mathbf{t}_1 \\ \mathbf{t}_2 \\ \vdots \\ \mathbf{t}_i \\ \vdots \\ \mathbf{t}_m \end{bmatrix} = \begin{bmatrix} 1 & t_1 \\ 1 & t_2 \\ \vdots & \vdots \\ 1 & t_k \\ \vdots & \vdots \\ 1 & t_m \end{bmatrix} \tag{18}$$

$$\mathbf{y} = [y_0 \quad y_1 \quad \dots \quad y_i \quad \dots \quad y_n]^T \tag{19}$$

so result is still the same

$$l_i(x; X) = \text{the } i^{\text{th}} \text{ column of } \mathbf{x}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \tag{20}$$

3 Polynomial Regression

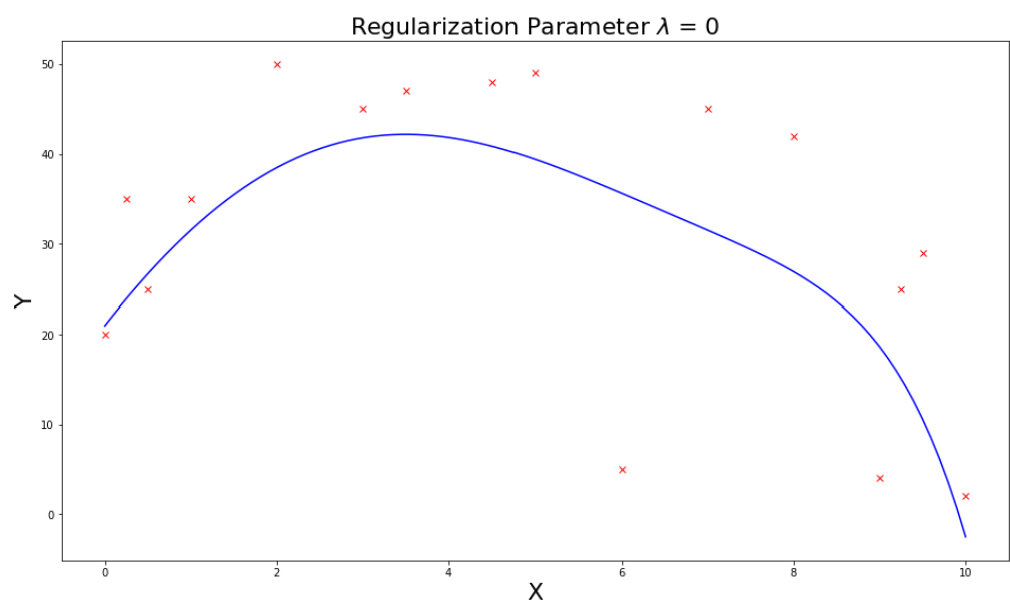


Figure 1: $\lambda = 0$

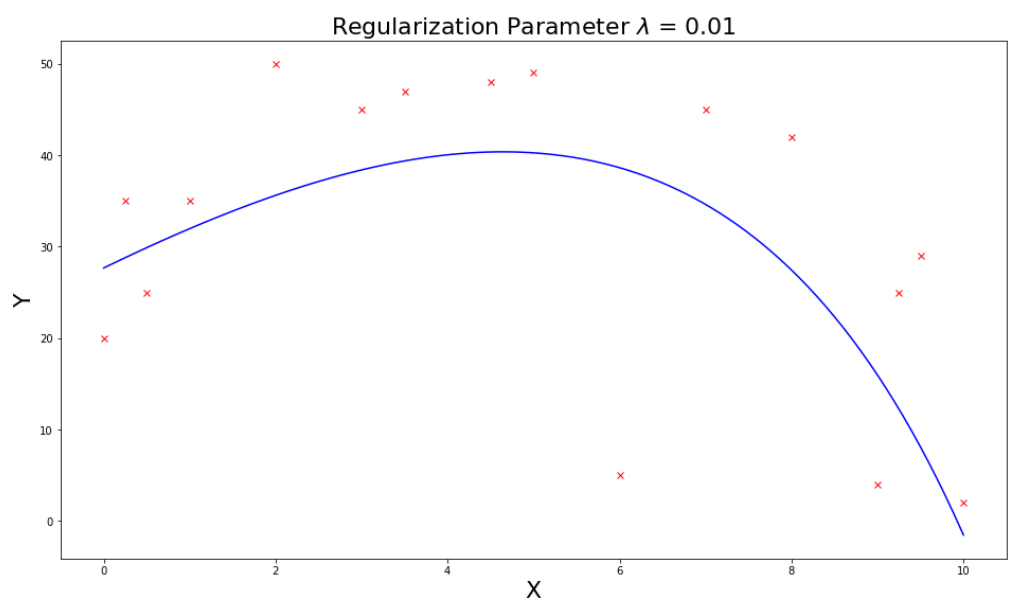


Figure 2: $\lambda = 0.01$