

1. Regularization is a technique to prevent overfitting by penalizing complex models. Weights that close to 0 is preferred, and it penalizes those ~~data~~ weights have high absolute value since $\lambda > 0$.

2. When $\lambda = 0$, it ignores to penalize w .
 When $\lambda \rightarrow +\infty$, it much prefer to penalize w .
 When $\lambda \rightarrow -\infty$, it prefer a w that much more fit the given training data.

$$3. \frac{\partial J}{\partial w_0} = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial w_0} (y_i - h w(x_i))^2 + \frac{\partial}{\partial w_0} \sum_{j=1}^p w_j^2$$

$$= \frac{1}{n} \sum_{i=1}^n -2(y_i - h w(x_i)) + 0$$

$$= -\frac{2}{n} \sum_{i=1}^n (y_i - h w(x_i))$$

$$\frac{\partial J}{\partial w_1} = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial w_1} (y_i - h w(x_i))^2 + \frac{\partial}{\partial w_1} \sum_{j=1}^p w_j^2$$

$$\frac{\partial J}{\partial w_1} = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial w_1} (y_i - h w(x_i))^2 + \frac{\partial}{\partial w_1} \sum_{j=1}^p w_j^2$$

$$= \frac{1}{n} \sum_{i=1}^n -2(y_i - h w(x_i))(-x_i) + \frac{\partial}{\partial w_1} \lambda \sum_{j=1}^p w_j^2$$

4. The different between this gradient descent and regularized linear gradient descent, ~~it~~ it has one more ~~ed~~

$$3. \frac{\partial J}{\partial w_0} = \frac{1}{n} \frac{\partial}{\partial w_0} \sum_{i=1}^n (y_i - h w(x_i))^2 + \frac{\partial}{\partial w_0} \lambda w_0^2$$

$$= \frac{1}{n} (-2) \sum_{i=1}^n (y_i - h w(x_i))$$

$$= -\frac{2}{n} \sum_{i=1}^n (y_i - h w(x_i))$$

$$\frac{\partial J}{\partial w_j} = \frac{1}{n} \sum_{i=1}^n (y_i - h w(x_i)) (-x_i) + 2\lambda \sum_{j=1}^P w_j$$

$$= -\frac{x_i}{n} \sum_{i=1}^n (y_i - h w(x_i)) + 2\lambda w_j = -\frac{1}{n} \sum_{i=1}^n (y_i - h w(x_i)) (-x_i) + 2\lambda \sum_{j=1}^P w_j$$

4. The difference between two gradient descent is

$$w_0 \leftarrow w_0 + \alpha \sum_{i=1}^n (y_i - h w(x_i))$$

$$w_j \leftarrow w_j + \alpha \left[\sum_{i=1}^n (y_i - h w(x_i)) \cdot x_i \right] + \alpha \left[2\lambda \sum_{j=1}^P w_j \right]$$

which has one more element considered in w_j

When ~~the~~ update w_0 , there is no difference,

but update w_j will become larger since $\lambda > 0$.

When w is positive, the loss function will get a large value
when w is negative, the loss function decreased which means
more accurate.

5. If designer want the agent be more accurate on given training dataset, it is better to set 0 for λ
If designer ~~don't~~ doesn't want overfitting, it is better to choose a λ larger than 0.

Perceptron

$$① \quad Z = h w_1 = w_1 x_1 + w_2 x_2 + w_3 x_3 + b x_0 = 0 \Rightarrow 1$$

$$w_1, w_2, w_3, b = 0$$

$$② \quad Z = h w_2 = w_1 x_1 + w_2 x_2 + w_3 x_3 + b x_0 = 0 \Rightarrow 1$$

$$w_1, w_2, w_3, b = 0$$

$$③ \quad Z = 0 \Rightarrow 1$$

$$w_1, w_2, w_3, b = 0$$

$$④ \quad Z = 0 \Rightarrow 1$$

$$w_1, w_2, w_3, b = 0$$

$$⑤ \quad Z = 0 \Rightarrow 1$$

$$w_1 = 0.5(0-1) \times 1 + 0 = -0.5 \quad w_2 = 0.5(0-1) \times 0 + 0 = 0$$

$$w_3 = 0.5(0-1) \times (-2) + 0 = 1 \quad b = 0.5(0-1) + 0 = -0.5$$

$$⑥ \quad Z = -0.5 \times (-1) + 0 + 1 + 0.5 = 2 \Rightarrow 1$$

$$w_1 = 0.5(0-1) \times 1 + -0.5 = 0 \quad w_2 = 0.5(0-1) \times -1 + 0 = 0.5$$

$$w_3 = 0.5(0-1) \times 1 + 0 = 0.5 \quad b = 0.5(0-1) + 0.5 = 0$$

$$⑦ \quad Z = 0 - 4 \times 0.5 - 1 = -3 \Rightarrow 0$$

$$w_1 = 0.5(0-0) \times 0 + 0.5 = 0.5 \quad w_2 = 0.5(0-0) \times -4 + 0.5 = 0.5$$

$$w_3 = 0.5(0-0) \times 0 + 0.5 = 0.5 \quad b = 0.5(0-0) - 1 = -1$$

$$⑧ \quad Z = 1 \times 0 + 0 \times 0.5 - 3 \times 0.5 - 1 = -2.5 \Rightarrow 0$$

$$w_1 = 0.5(0-0) \times 1 + 0 = 0 \quad w_2 = 0.5(0-0) \times 0 + 0.5 = 0.5$$

$$w_3 = 0.5(0-0) \times -3 + 0.5 = 0.5 \quad b = 0.5(0-0) - 1 = -1$$

$$⑨ \quad Z = 4 \times 0.5 + 0.5 - 1 = 1.5 \Rightarrow 1$$

$$w_1 = 0 \quad w_2 = 0.5 \quad w_3 = 0.5 \quad b = -1$$

$$⑩ \quad Z = 2 \times 0.5 + 0.5 \times 3 - 1 = 1.5 \Rightarrow 1$$

$$w_1 = 0 \quad w_2 = 0.5 \quad w_3 = 0.5 \quad b = -1$$

$$⑪ \quad Z = 0.5 \times 1 = 0.5 \Rightarrow 0$$

$$w_1 = 0.5(1-0) \times 0 + 0 = 0 \quad w_2 = 0.5(1-0) \times 0 + 0.5 = 0.5$$

$$w_3 = 0.5(1-0) \times 1 + 0.5 = 1 \quad b = 0.5(1-0) - 1 = -0.5$$

$$⑫ \quad Z = 0 \times 1 + 0 \times 0.5 - 3 \times 1 - 0.5 = -3.5 \Rightarrow 0$$

$$w_1 = 0.5 \quad w_2 = 0.5 \quad w_3 = 0.5$$

$$⑬ \quad Z = 0 - 1 \times 0 + 4 \times 0.5 + 0 \times 1 - 0.5 = 1.5 \Rightarrow 1$$

$$w_1 = 0 \quad w_2 = 0.5 \quad w_3 = 1 \quad b = -0.5$$

$$\textcircled{5} Z = 1 \times 0 + 0 \times 0.5 - 2 \times 1 - 0.5 = -2.5 \Rightarrow 0$$

$$w_1 = 0 \quad w_2 = 0.5 \quad w_3 = 1 \quad b = -0.5$$

$$\textcircled{6} Z = -1 \times 0 + 4 \times 0.5 - 1 \times -0.5 = \overset{-2}{-0.5} = 0$$

$$w_1 = 0 \quad w_2 = 0.5 \quad w_3 = 1 \quad b = -0.5$$

$$\textcircled{7} Z = 0 \times 0 - 4 \times 0.5 - 0.5 = -2.5 \Rightarrow 0$$

$$w_1 = 0 \quad w_2 = 0.5 \quad w_3 = 1 \quad b = -0.5$$

$$\textcircled{8} Z = 1 \times 0 - 0 \times 0.5 - 3 \times -0.5 = -3.5 \Rightarrow 0$$

$$w_1 = 0 \quad w_2 = 0.5 \quad w_3 = 1 \quad b = -0.5$$

3.

Sample	N_1	N_2	N_3	N_4	N_5	N_6	N_7	Output
x_1	0	1	1	0	1	1	0	0
x_2	0	1	1	1	0	0	0	0
x_3	1	0	1	1	0	1	0	1
x_4	1	0	0	0	0	1	1	0

where ~~there~~ N_i is i th cell, totally 7 cells.
 empty cell (cell with white color) is 0,
 colored cell (cell with gray color) is 1.

~~total~~ one gray area is defined by all
 connected cells colored. If the ~~total~~ number of
 a ~~cell~~ gray areas are more than
 2, the output is 1, otherwise, 0.

PS: N_1 and N_7 will sign as connected.

learning rate $\alpha = 0.5$, starts ~~in~~ with initial weights of 0,
 bias neuron $x_0 = 1$