Since a new regularization term was introduced, the derivative has been changed.

$$\frac{\partial J(\vec{w},b)}{\partial w_j} \longrightarrow \frac{1}{m} \sum_{i=1}^{m} (f_{\vec{w},b}(x^{(i)}) - y^{(i)}) \, x_j^{(i)} + \frac{\lambda}{m} w_j$$

$$\frac{\partial J(\vec{w},b)}{\partial b} \longrightarrow \frac{1}{m} \sum_{i=1}^{m} f_{\vec{w},b}(x^{(i)}) - y^{(i)}$$

it remains the same because we don't regularize b

Gradient Descent :-

repeat {

$$w_j = w_j - \alpha \boxed{\frac{\partial J(\vec{w},b)}{\partial w_j}} \longrightarrow \frac{1}{m} \sum_{i=1}^{m} (f_{\vec{w},b}(x^{(i)}) - y^{(i)}) \, x_j^{(i)} + \frac{\lambda}{m} w_j$$

$$b = b - \alpha \boxed{\frac{\partial J(\vec{w},b)}{\partial b}} \longrightarrow \frac{1}{m} \sum_{i=1}^{m} f_{\vec{w},b}(x^{(i)}) - y^{(i)}$$

} simultaneous update

How is regularization shrinking all the parameters?

Gradient descent is as follows:-

$$w_j = w_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} (f_{\vec{w},b}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \sum_{j=1}^{n} w_j \right]$$

$$\underline{\frac{\partial J(\vec{w},b)}{\partial w_j}}$$

$$\Rightarrow w_j = w_j - \frac{\alpha}{m} \sum_{i=1}^{m} (f_{\vec{w},b}(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\alpha \lambda}{m} w_j$$

On rearrangement

$$\Rightarrow w_j = w_j - \frac{\alpha \lambda}{m} w_j - \frac{1}{m} \sum_{i=1}^{m} (f_{\vec{w},b}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\Rightarrow w_j = w_j \left[ 1 - \alpha \frac{\lambda}{m} \right] - \frac{\alpha}{m} \sum_{i=1}^{m} (f_{\vec{w},b}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

taking $w_j$ common            usual term

Let's take $\alpha = 0.01$, $\lambda = 1$ and $m = 50$, for example.

$$w_j = w_j \left[ 1 - 0.0002 \right] - \frac{\alpha}{m} \sum_{i=1}^{m} (f_{\vec{w},b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

$$w_j = w_j \underbrace{(0.9998)}_{\text{very small}} - \frac{\alpha}{m} \sum_{i=1}^{m} (f_{\vec{w},b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

\# Earlier we used to subtract just the usual term $\left(\dfrac{\alpha}{m} \sum\limits_{i=1}^{m} (f_{\vec{w},b}(x^{(i)}) - y^{(i)}) x^{(i)}\right)$ from $w_j$, but now

we have to multiply a small number like $0.9998$ ($<1$) to $w_j$ and then subtract the usual term.

This makes the $w_j$ shrinking process more effective and faster by decreasing $w_j$ in each iteration by a little bit.

# How we get the derivative term (optional)

$$\frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{d}{dw_j}\left[\frac{1}{2m}\sum_{i=1}^{m} \underbrace{\left(f(\vec{x}^{(i)}) - y^{(i)}\right)^2}_{\vec{w}\cdot\vec{x}^{(i)}+b} + \frac{\lambda}{2m}\sum_{j=1}^{n} w_j^2\right]$$

$$= \frac{1}{2m}\sum_{i=1}^{m}\left[(\vec{w}\cdot\vec{x}^{(i)} + b - y^{(i)})\,2x_j^{(i)}\right] + \frac{\lambda}{2m}\,2w_j \qquad \text{No } \sum_{j=1}^{n}$$

$$= \frac{1}{m}\sum_{i=1}^{m}\left[(\underbrace{\vec{w}\cdot\vec{x}^{(i)} + b}_{f(\vec{x})} - y^{(i)})\,x_j^{(i)}\right] + \frac{\lambda}{m}\,w_j$$

no need for summation because in gradient descent $w$ will automatically be updated.