1-) For every ML problem we should define

    1-) Hypothesis = Linear hypothesis with non linear parameters.

    2-) Loss = Hinge loss which can not be differentiable at some point

    3-) Optimize the loss

Since differentiation of hingle loss has discontinuity therefore there is no closed form solution. We should use sub gradient.

$$\text{Hypothesis} = h_\theta(x) = W^T \cdot X \qquad \text{where} \quad X = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_6^{(i)} \\ 1 \end{bmatrix}$$

Note: $x^{(i)}$ is the i'th instance of data points

$$\text{loss} \quad \ell(\theta) = \max(1 - y h_\theta(x), 0)$$

Our optimization

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^{n} \max\left(1 - y^{(i)} w^T x^{(i)}\right) \Rightarrow J(w) = \frac{1}{n} \sum_{i=1}^{n} \max\left(1 - y^{(i)} w^T x^{(i)}, 0\right)$$

$$\nabla_w \ell\left(y^{(i)} w^T x^{(i)}\right) = \ell'\left(y^{(i)} w^T x^{(i)}\right) y^{(i)} x^{(i)} = \left( \begin{cases} 0, & y^{(i)} w^T x^{(i)} > 1 \\ -1, & y^{(i)} w^T x^{(i)} < 1 \\ \text{undefined}, & y^{(i)} w^T x^{(i)} = 1 \end{cases} \right) y^{(i)} x^{(i)}$$

$$= \begin{cases} 0, & y^{(i)} w^T x^{(i)} > 1 \\ -y^{(i)} x^{(i)}, & y^{(i)} w^T x^{(i)} < 1 \\ \text{undefined}, & y^{(i)} w^T x^{(i)} = 1 \end{cases} = \nabla_w \ell\left(y^{(i)} w^T x^{(i)}\right)$$

$$\nabla_w J(w) = \nabla_w \left( \frac{1}{n} \sum_{i=1}^{n} \ell\left(y^{(i)} w^T x^{(i)}\right) \right) = \frac{1}{n} \sum_{i=1}^{n} \nabla_w \ell\left(y_i w^T x_i\right)$$

$$= \begin{cases} \frac{1}{n}\left( \sum_{i=\, y^{(i)} w^T x < 1} -y^{(i)} x^{(i)} \right), & \text{all } y^{(i)} w^T x^{(i)} \neq 1 \\ \text{undefined}, & \text{otherwise} \end{cases}$$

Pseudo code:
```
w ← <0,0,----0>   ← weights
b ← 0             ← bias
for iter = 1 --- epochs do
    g ← <0,0----0>  ← gradients
    g_b ← 0
    for all (x,y) do
        if y(wx + b) ≤ 1
            g ← g + yx
            g_b ← g + y
    end for
    w ← w + λg      → step size
    b ← b + λg_b
```

**2-)**

Output of the $z_4 = \sigma(z_4) = \frac{1}{1+e^{-z_4}}$

**2.1-)** $L = -y\log(\hat{y}) - (1-y)\log(1-\hat{y})$ $\quad \frac{\partial L}{\partial \hat{y}} = \frac{-y}{\hat{y}} + \frac{y-1}{1-\hat{y}} = \frac{\hat{y}-y}{\hat{y}(1-\hat{y})}$

$\hat{y} = \sigma(z_4) = \frac{1}{1+e^{-z_4}}$ $\quad \frac{\partial \hat{y}}{\partial z_4} = \sigma(z_4)\cdot(1-\sigma(z_4))$

So $\frac{\partial L}{\partial \sigma(z_4)} = \frac{\sigma(z_4)(1-\sigma(z_4)) - y}{\sigma(z_4)(1-\sigma(z_4))\,(1-\sigma(z_4)(1-\sigma(z_4))}$

**2.2-)** $\frac{\partial L}{\partial W_{14}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_4} \cdot \frac{\partial z_4}{\partial W_{14}} = \delta z_4 \cdot \frac{\partial z_4}{\partial W_{14}} = \delta z_4 \cdot a_1 = \frac{\partial L}{\partial W_{14}}$

Given $\delta z_4 = \frac{\partial L}{\partial z_4}$ $\qquad\qquad a_1$

**2.3-)** $\frac{\partial L}{\partial W_{11}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_4} \cdot \frac{\partial z_4}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial W_{11}} = \delta a_1 \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial W_{11}}$

Given $\delta a_1 = \frac{\partial L}{\partial a_1}$

$a_1 = ReLU(z_1) = \max(0, z_1) = \begin{cases} 0 & \text{if } z_1 < 0 \\ z_1 & \text{if } z_1 > 0 \end{cases}$

$\frac{\partial a_1}{\partial z_1} = \begin{cases} 0 & \text{if } z_1 < 0 \\ 1 & \text{if } z_1 > 0 \end{cases}$

$= \delta a_1 \cdot \begin{cases} 0 & \text{if } z_1 < 0 \\ 1 & \text{if } z_1 > 0 \end{cases} \cdot x_1$

$= \begin{cases} 0 & \text{if } z_1 < 0 \\ \delta a_1 \cdot x_1 & \text{if } z_1 > 0 \end{cases}$

$z_1 = W_{11} \cdot x_1 + b_1$

$\frac{\partial z_1}{\partial W_{11}} = x_1$

2.4-) L2 regularization is only applied on weights in general. Biases are omitted.

We add L2 norm on cost function with $\lambda$ regularization parameter.

First gradient: Our new cost $L = -y \log(\hat{y}) - (1-y) \log(1-\hat{y}) + \dfrac{\lambda}{2} \sum_w w^2$

So for every partial derivate on weight becomes $\boxed{\dfrac{\partial L}{\partial W_{ij}} = \dfrac{\partial L}{\partial W_{ij}} + \lambda \cdot w}$

Since we don't use any $W$ in $\dfrac{\partial L}{\partial \sigma(z_4)}$ there is no change.

Note: I assume same upstream gradients for the following gradients.

Second gradient: $\dfrac{\partial L}{\partial w_{14}} = \underbrace{\dfrac{\partial L}{\partial \hat{y}} \cdot \dfrac{\partial \hat{y}}{\partial z_4}}_{\delta z_4} \cdot \underbrace{\dfrac{\partial z_4}{\partial w_{14}}}_{a_1} = \delta z_4 \cdot \left(a_1 + \lambda \cdot w_{14}\right) = \boxed{\delta z_4 \cdot a_1 + \delta z_4 \cdot \lambda \cdot W_{14}}$

$\delta z_4 = \dfrac{\partial L}{\partial z_4}$

Third gradient: We have found that $\dfrac{\partial L}{\partial w_{11}} = \delta a_1 \cdot \dfrac{\partial a_1}{\partial z_1} \cdot \dfrac{\partial z_1}{\partial w_{11}}$

$\dfrac{\partial a_1}{\partial z_1} = \begin{cases} 0 & \text{if } z_1 < 0 \\ 1 & \text{if } z_1 > 0 \end{cases}$

$z_1 = w_{11} x_1 + b$

$\dfrac{\partial z_1}{\partial w_{11}} = x_1 + \lambda \cdot w_{11}$

$= \delta a_1 \cdot \begin{cases} 0 & \text{if } z_1 < 0 \\ 1 & \text{if } z_1 > 0 \end{cases} \cdot (x_1 + \lambda w_{11})$

$= \begin{cases} 0 & \text{if } z_1 < 0 \\ \delta a_1 (x_1 + \lambda w_{11}) & \text{if } z_1 > 0 \end{cases}$

So in general regularization term regularize the derivates with $W_{ij}$ weights.

So it changes when $w$ is involved in the gradient.