

# 1 Source

To the best of my recollection, the details of the logistic regression update rule math were not reviewed in class. As such, I watched the lectures from [Andrew Ng's deep learning class](#). Here is what I *believe* the update rule should be. If there is an error in the logic, please let me know.

## 2 Glossary of Notation

- $\mathbf{w}_t$  – Weight vector for epoch  $t$
- $J(\mathbf{w}, \mathbf{x})$  – Cost function
- $\beta$  – Learning rate
- $\mathcal{L}$  – Loss function
- $\hat{y}$  – Predicted output value
- $y$  – Expected (target) classification value
- $\sigma(z)$  – Sigmoid function  $\left(\frac{1}{1+e^{-z}}\right)$  with respect to  $z$  (i.e.,  $\mathbf{w}^T \mathbf{x}$ ).

## 3 $w$ Update Rules with Squared Loss

My understanding of the *batch* update rule is shown in Eq. (1).

$$\mathbf{w}_{t+1} := \mathbf{w}_t - \beta \cdot \frac{\partial J(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}} \quad (1)$$

The learning rate is defined as:

$$\beta := \eta \cdot t^{-\alpha}$$

where  $\alpha = 0.9$ . The cost function is the average loss as shown in Eq. (2).

$$J(\mathbf{w}, \mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\mathbf{w}, \mathbf{x}) \quad (2)$$

The loss function is the squared loss and uses the same regularizer as in homework #1 as shown in Eq. (3).

$$\mathcal{L}(\hat{y}, y, \mathbf{w}) = \frac{1}{2}(\hat{y} - y)^2 + \lambda \|\mathbf{w}\| \quad (3)$$

The predicted value  $\hat{y}$  is

$$\hat{y} = \sigma(\mathbf{w}^T \mathbf{x}).$$

The derivative of the loss function  $\mathcal{L}$  is:

$$\frac{\partial \mathcal{L}(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}} = (\hat{y} - y) \frac{\partial \hat{y}}{\partial \mathbf{w}} + \lambda \mathbf{w}. \quad (4)$$

Via the chain rule, we solve the derivative of the sigmoid function:

$$\frac{\partial \hat{y}(z)}{\partial \mathbf{w}} = \sigma(z)(1 - \sigma(z)) = \frac{e^{-z}}{(1 + e^{-z})^2} \frac{\partial z}{\partial \mathbf{w}} \quad (5)$$

Applying the chain rule again yields:

$$\frac{\partial z}{\partial \mathbf{w}} = \mathbf{x} \quad (6)$$

Combining Eq. (4), (5), and (6) shows the complete derivative of the loss function.

$$\frac{\partial \mathcal{L}(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}} = \left( \sigma(\mathbf{w}^\top \mathbf{x}) - y \right) \left( \frac{e^{-z}}{(1 + e^{-z})^2} \right) \mathbf{x} + \lambda \mathbf{w} \quad (7)$$

$$= \left( \sigma(\mathbf{w}^\top \mathbf{x}) - y \right) \left( \frac{e^{-\mathbf{w}^\top \mathbf{x}}}{(1 + e^{-\mathbf{w}^\top \mathbf{x}})^2} \right) \mathbf{x} + \lambda \mathbf{w} \quad (8)$$

For mathematical simplicity, the identity for  $\sigma'(z)$  allows for a simpler form in Eq. (9).<sup>1</sup>

$$\frac{\partial \mathcal{L}(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}} = \left( \sigma(\mathbf{w}^\top \mathbf{x}) - y \right) \sigma(\mathbf{w}^\top \mathbf{x}) \left( 1 - \sigma(\mathbf{w}^\top \mathbf{x}) \right) \mathbf{x} + \lambda \mathbf{w} \quad (9)$$

## 4 $w$ Update Rules with Logistic Loss

The more common loss function I see for logistic regression is in Eq. (10).

$$\mathcal{L}(\hat{y}, y, \mathbf{w}) = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})) \quad (10)$$

The derivative of the logistic loss function is shown below in Eq. (11).

$$\frac{\partial \mathcal{L}(\hat{y}, y, \mathbf{w})}{\partial \mathbf{w}} = - \left( \frac{y}{\hat{y}} - \frac{1 - y}{1 - \hat{y}} \right) \frac{\partial \hat{y}}{\partial \mathbf{w}} \quad (11)$$

We know the derivative of  $\hat{y}$  from Eq. (5). Substituting that we get:

$$\frac{\partial \mathcal{L}(\hat{y}, y, \mathbf{w})}{\partial \mathbf{w}} = - \left( \frac{y}{\hat{y}} - \frac{1 - y}{1 - \hat{y}} \right) \hat{y}(1 - \hat{y}) \frac{\partial \mathbf{z}}{\partial \mathbf{w}} \quad (12)$$

$$= (-y(1 - \hat{y}) + (1 - y)\hat{y}) \frac{\partial \mathbf{z}}{\partial \mathbf{w}} \quad (13)$$

$$= (\hat{y} - y) \frac{\partial \mathbf{z}}{\partial \mathbf{w}}. \quad (14)$$

The complete derivative then is in Eq. (15).

$$\frac{\partial \mathcal{L}(\hat{y}, y, \mathbf{w})}{\partial \mathbf{w}} = (\hat{y} - y) \mathbf{x}. \quad (15)$$

---

<sup>1</sup>I am not considering the transposes. That math I would need to think more about.