

Logistic Regression Gradient Descent Detailed Derivative Derivations

1. Variables:

The following table provides an inventory of all the variables used:

#	Name	Description
1	m	Sample size
2	n_x	Number of input features
3	$x_j^{(i)}$	j th input feature of the i th sample
4	$\vec{x}^{(i)}$	Equal to $[x_1^{(i)}, x_2^{(i)}, \dots, x_{n_x}^{(i)}]^T$
5	X	Equal to $[\vec{x}^{(1)}, \vec{x}^{(2)}, \dots, \vec{x}^{(m)}]$
6	w_j	j th input feature weight
7	\vec{w}	Equal to $[w_1, w_2, \dots, w_{n_x}]^T$
8	b	Bias
9	$\mathbf{1}_{n \times p}$	Matrix of ones with n rows and p columns
10	$z^{(i)}$	Equal to $\vec{w}^T \vec{x}^{(i)} + b \mathbf{1}_{m \times 1}$
11	\vec{z}	Equal to $[z^{(1)}, z^{(2)}, \dots, z^{(m)}]$
12	$g(x)$	Transformation function on variable x . Here, we set $g(x) = \frac{1}{1+e^{-x}}$, which is the sigmoid function
13	$a^{(i)}$	Activation value of the i th sample. Equal to $g(z^{(i)})$
14	$g^*(\vec{v}_{1 \times p})$	For $\vec{v}_{1 \times p} = [x_1, x_2, \dots, x_p]$, equal to $[g(x_1), g(x_2), \dots, g(x_p)]$
15	\vec{a}	Equal to $g^*(\vec{z})$
16	$y^{(i)}$	Output variable of the i th sample. Possible values $\in \{0,1\}$
17	Y	Equal to $[y^{(1)}, y^{(2)}, \dots, y^{(m)}]$
18	$\mathcal{L}(a^{(i)}, y^{(i)})$	Loss value of the i th sample. Set equal to $-(y^{(i)} \ln(a^{(i)}) + (1 - y^{(i)}) \ln(1 - a^{(i)}))$
19	\mathcal{J}	Loss value of the entire training dataset. Set equal to $\frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^{(i)}, y^{(i)})$

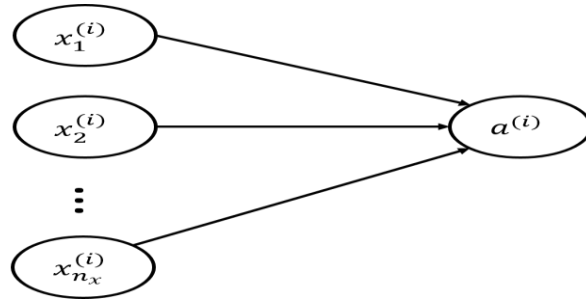
2. Vector Derivative Convention

We use the so-called numerator layout convention: let $\vec{v}_{n \times 1} = [v_1, v_2, \dots, v_n]^T$, $\vec{u}_{1 \times p} = [u_1, u_2, \dots, u_p]$, $f_u(\vec{u}_{1 \times p})$ a function of the vector $\vec{u}_{1 \times p}$ mapping onto \mathbb{R} . Furthermore, let $\vec{v}_{n \times 1} = h(\vec{u}_{1 \times p})$. We thus have the following results based on the chosen convention:

- $$\frac{\partial h(\bar{u}_{1 \times p})}{\partial u_i} = \frac{\partial \bar{v}_{n \times 1}}{\partial u_i} = \begin{bmatrix} \frac{\partial v_1}{\partial u_i} \\ \frac{\partial v_n}{\partial u_i} \\ \vdots \\ \frac{\partial v_n}{\partial u_i} \end{bmatrix}$$
- $$\frac{\partial (h(\bar{u}_{1 \times p}))^T}{\partial u_i} = \frac{\partial (\bar{v}_{n \times 1})^T}{\partial u_i} = \begin{bmatrix} \frac{\partial v_1}{\partial u_i} \\ \frac{\partial v_n}{\partial u_i} \\ \vdots \\ \frac{\partial v_n}{\partial u_i} \end{bmatrix}^T$$
- $$\frac{\partial f_u(\bar{u}_{1 \times p})}{\partial \bar{u}_{1 \times p}} = \begin{bmatrix} \frac{\partial f_u(\bar{u}_{1 \times p})}{\partial u_1} & \frac{\partial f_u(\bar{u}_{1 \times p})}{\partial u_2} & \dots & \frac{\partial f_u(\bar{u}_{1 \times p})}{\partial u_p} \end{bmatrix}$$
- $$\frac{\partial f_u(\bar{u}_{1 \times p})}{\partial (\bar{u}_{1 \times p})^T} = \begin{bmatrix} \frac{\partial f_u(\bar{u}_{1 \times p})}{\partial u_1} & \frac{\partial f_u(\bar{u}_{1 \times p})}{\partial u_2} & \dots & \frac{\partial f_u(\bar{u}_{1 \times p})}{\partial u_p} \end{bmatrix}^T$$
- $$\frac{\partial h(\bar{u}_{1 \times p})}{\partial \bar{u}_{1 \times p}} = \begin{bmatrix} \frac{\partial h(\bar{u}_{1 \times p})}{\partial u_1} & \frac{\partial h(\bar{u}_{1 \times p})}{\partial u_2} & \dots & \frac{\partial h(\bar{u}_{1 \times p})}{\partial u_p} \end{bmatrix} = \begin{bmatrix} \frac{\partial v_1}{\partial u_1} & \frac{\partial v_1}{\partial u_2} & \dots & \frac{\partial v_1}{\partial u_p} \\ \frac{\partial v_2}{\partial u_1} & \frac{\partial v_2}{\partial u_2} & \dots & \frac{\partial v_2}{\partial u_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial v_n}{\partial u_1} & \frac{\partial v_n}{\partial u_2} & \dots & \frac{\partial v_n}{\partial u_p} \end{bmatrix}$$
- $$\frac{\partial h(\bar{u}_{1 \times p})^T}{\partial \bar{u}_{1 \times p}^T} = \begin{bmatrix} \frac{\partial h(\bar{u}_{1 \times p})}{\partial u_1} & \frac{\partial h(\bar{u}_{1 \times p})}{\partial u_2} & \dots & \frac{\partial h(\bar{u}_{1 \times p})}{\partial u_p} \end{bmatrix}^T = \begin{bmatrix} \frac{\partial v_1}{\partial u_1} & \frac{\partial v_2}{\partial u_1} & \dots & \frac{\partial v_n}{\partial u_1} \\ \frac{\partial v_1}{\partial u_2} & \frac{\partial v_2}{\partial u_2} & \dots & \frac{\partial v_n}{\partial u_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial v_1}{\partial u_p} & \frac{\partial v_2}{\partial u_p} & \dots & \frac{\partial v_n}{\partial u_p} \end{bmatrix} =$$
- $$\begin{bmatrix} \frac{\partial v_1}{\partial u_1} & \frac{\partial v_1}{\partial u_2} & \dots & \frac{\partial v_1}{\partial u_p} \\ \frac{\partial v_2}{\partial u_1} & \frac{\partial v_2}{\partial u_2} & \dots & \frac{\partial v_2}{\partial u_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial v_n}{\partial u_1} & \frac{\partial v_n}{\partial u_2} & \dots & \frac{\partial v_n}{\partial u_p} \end{bmatrix}^T$$

3. Model Setup

The following diagram provides a visual representation of the model for training sample i :



The following formulas describe the model behaviour:

- Individual sample i :
 - $a^{(i)} = g(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}}$
 - $z^{(i)} = \vec{w}^T \vec{x}^{(i)} + b$
 - $\mathcal{L}(a^{(i)}, y^{(i)}) = -(y^{(i)} \ln(a^{(i)}) + (1 - y^{(i)}) \ln(1 - a^{(i)}))$
- Training dataset:
 - $\vec{a} = [a^{(1)}, a^{(2)}, \dots, a^{(m)}] = [g(z^{(1)}), g(z^{(2)}), \dots, g(z^{(m)})] = g^*(\vec{z})$
 - $\vec{z} = [z^{(1)}, z^{(2)}, \dots, z^{(m)}] = [\vec{w}^T \vec{x}^{(1)} + b, \vec{w}^T \vec{x}^{(2)} + b, \dots, \vec{w}^T \vec{x}^{(m)} + b]$
 - $\mathcal{J} = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^{(i)}, y^{(i)})$

4. Goal

Find the parameter values $w_1, w_2, \dots, w_{n_x}, b$ that minimize \mathcal{J}

5. Method

Gradient descent: An iterative algorithm that recalibrates the parameter values systematically reducing \mathcal{J} . Formulas (square bracket superscript reflects iteration number):

- w_j 's, where $j = 1, 2, \dots, n_x$:
 - $w_j^{[0]} = \text{random number (typically chosen to be close to zero)} \rightarrow \text{Vector form: } \vec{w}^{[0]}$
 - $w_j^{[k+1]} = w_j^{[k]} - \alpha \times \frac{\partial \mathcal{J}^{[k]}}{\partial w_j^{[k]}}$, where α is a set learning rate (e.g., 0.005), and $k = 0, 1, \dots, N$, where N is the total number of iterations $\rightarrow \text{Vector form: } \vec{w}^{[k+1]} = \vec{w}^{[k]} - \alpha \times \frac{\partial \mathcal{J}^{[k]}}{\partial \vec{w}^{[k]}}$
- b :
 - $b^{[0]} = \text{random number (typically chosen to be zero)}$
 - $b^{[k+1]} = b^{[k]} - \alpha \times \frac{\partial \mathcal{J}^{[k]}}{\partial b^{[k]}}$, where α is a set learning rate (e.g., 0.005), and $k = 0, 1, \dots, N$, where N is the total number of iterations

Notes:

- $\mathcal{J}^{[k]} = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^{(i)[k]}, y^{(i)})$; $y^{(i)}$ doesn't have a superscript k since its value does not change as we iterate

- We will prove that $\mathcal{J}^{[k]}$ is convex with respect to w_j 's (making one minor assumption to guarantee strict convexity) and b (strictly convex), justifying putting a minus sign next to the α parameters

5.1 Deriving $\frac{\partial \mathcal{J}^{[k]}}{\partial w_j^{[k]}}$ and $\frac{\partial \mathcal{J}^{[k]}}{\partial b^{[k]}}$

Via the chain rule, we have the following results:

- $\frac{\partial \mathcal{J}^{[k]}}{\partial w_j^{[k]}} = \frac{\partial \mathcal{J}^{[k]}}{\partial \vec{a}^{[k]}} \times \frac{\partial \vec{a}^{[k]}}{\partial \vec{z}^{[k]T}} \times \frac{\partial \vec{z}^{[k]T}}{\partial w_j^{[k]}}$ for $j = 1, 2, \dots, n_x$ and $k = 0, 1, \dots, N$
- $\frac{\partial \mathcal{J}^{[k]}}{\partial b^{[k]}} = \frac{\partial \mathcal{J}^{[k]}}{\partial \vec{a}^{[k]}} \times \frac{\partial \vec{a}^{[k]}}{\partial \vec{z}^{[k]T}} \times \frac{\partial \vec{z}^{[k]T}}{\partial b^{[k]}}$ for $k = 0, 1, \dots, N$

We now derive the mathematical expression of each component separately:

5.1.1 Deriving $\frac{\partial \mathcal{J}^{[k]}}{\partial \vec{a}^{[k]}}$

$$\begin{aligned} \frac{\partial \mathcal{J}^{[k]}}{\partial \vec{a}^{[k]T}} &= \left[\frac{\partial \mathcal{J}^{[k]}}{\partial a^{(1)[k]}}, \frac{\partial \mathcal{J}^{[k]}}{\partial a^{(2)[k]}}, \dots, \frac{\partial \mathcal{J}^{[k]}}{\partial a^{(m)[k]}} \right] \\ &= \left[\frac{\partial \frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^{(i)[k]}, y^{(i)})}{\partial a^{(1)[k]}}, \frac{\partial \frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^{(i)[k]}, y^{(i)})}{\partial a^{(2)[k]}}, \dots, \frac{\partial \frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^{(i)[k]}, y^{(i)})}{\partial a^{(m)[k]}} \right] \\ &= \frac{1}{m} \left[\frac{\partial \mathcal{L}(a^{(1)[k]}, y^{(1)})}{\partial a^{(1)[k]}}, \frac{\partial \mathcal{L}(a^{(2)[k]}, y^{(2)})}{\partial a^{(2)[k]}}, \dots, \frac{\partial \mathcal{L}(a^{(m)[k]}, y^{(m)})}{\partial a^{(m)[k]}} \right] \end{aligned}$$

For $i = 1, 2, \dots, m$,

$$\begin{aligned} \frac{\partial \mathcal{L}(a^{(i)[k]}, y^{(i)})}{\partial a^{(i)[k]}} &= - \left(\frac{y^{(i)}}{a^{(i)[k]}} - \frac{1 - y^{(i)}}{1 - a^{(i)[k]}} \right) \\ &= \frac{-(y^{(i)} \times (1 - a^{(i)[k]}) - a^{(i)[k]} \times (1 - y^{(i)}))}{a^{(i)[k]} \times (1 - a^{(i)[k]})} \\ &= \frac{-(y^{(i)} - y^{(i)} a^{(i)[k]} - a^{(i)[k]} + a^{(i)[k]} y^{(i)})}{a^{(i)[k]} \times (1 - a^{(i)[k]})} \\ &= \frac{a^{(i)[k]} - y^{(i)}}{a^{(i)[k]} \times (1 - a^{(i)[k]})} \end{aligned}$$

$$\text{Thus, } \frac{\partial \mathcal{J}^{[k]}}{\partial \vec{a}^{[k]}} = \frac{1}{m} \left[\frac{a^{(1)[k]} - y^{(1)}}{a^{(1)[k]} \times (1 - a^{(1)[k]})}, \frac{a^{(2)[k]} - y^{(2)}}{a^{(2)[k]} \times (1 - a^{(2)[k]})}, \dots, \frac{a^{(m)[k]} - y^{(m)}}{a^{(m)[k]} \times (1 - a^{(m)[k]})} \right]$$

5.1.2 Deriving $\frac{\partial \vec{a}^{[k]T}}{\partial \vec{z}^{[k]}}$

$$\frac{\partial \vec{a}^{[k]}}{\partial \vec{z}^{[k]T}} = \begin{bmatrix} \frac{\partial \vec{a}^{[k]}}{\partial z^{(1)[k]}} \\ \frac{\partial \vec{a}^{[k]}}{\partial z^{(2)[k]}} \\ \vdots \\ \frac{\partial \vec{a}^{[k]}}{\partial z^{(m)[k]}} \end{bmatrix} = \begin{bmatrix} \frac{\partial a^{(1)[k]}}{\partial z^{(1)[k]}} & \frac{\partial a^{(2)[k]}}{\partial z^{(1)[k]}} & \cdots & \frac{\partial a^{(m)[k]}}{\partial z^{(1)[k]}} \\ \frac{\partial a^{(1)[k]}}{\partial z^{(2)[k]}} & \frac{\partial a^{(2)[k]}}{\partial z^{(2)[k]}} & \cdots & \frac{\partial a^{(m)[k]}}{\partial z^{(2)[k]}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial a^{(1)[k]}}{\partial z^{(m)[k]}} & \frac{\partial a^{(2)[k]}}{\partial z^{(m)[k]}} & \cdots & \frac{\partial a^{(m)[k]}}{\partial z^{(m)[k]}} \end{bmatrix}$$

For $i = 1, 2, \dots, m$, and $j = 1, 2, \dots, m$, we have the following expression:

$$\frac{\partial a^{(i)[k]}}{\partial z^{(j)[k]}} = \frac{\partial \left(\frac{1}{1 + e^{-z^{(i)[k]}}} \right)}{\partial z^{(j)[k]}} = \begin{cases} \frac{e^{-z^{(i)[k]}}}{(1 + e^{-z^{(i)[k]}})^2} \text{ if } i = j \\ 0 \text{ if } i \neq j \end{cases}$$

For $i = j$, $\frac{e^{-z^{(i)[k]}}}{(1 + e^{-z^{(i)[k]}})^2} = \frac{1}{1 + e^{-z^{(i)[k]}}} \times \frac{e^{-z^{(i)[k]}}}{1 + e^{-z^{(i)[k]}}} = \left(\frac{1}{1 + e^{-z^{(i)[k]}}} \right) \times \left(1 - \frac{1}{1 + e^{-z^{(i)[k]}}} \right)$, and thus, we can express the expression in terms of $a^{(i)[k]}$:

$$\frac{\partial a^{(i)[k]}}{\partial z^{(j)[k]}} = \begin{cases} a^{(i)[k]} \times (1 - a^{(i)[k]}) \text{ if } i = j \\ 0 \text{ if } i \neq j \end{cases}$$

$$\text{Thus, } \frac{\partial \vec{a}^{[k]}}{\partial \vec{z}^{[k]T}} = \begin{bmatrix} a^{(1)[k]} \times (1 - a^{(1)[k]}) & 0 & \cdots & 0 \\ 0 & a^{(2)[k]} \times (1 - a^{(2)[k]}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a^{(m)[k]} \times (1 - a^{(m)[k]}) \end{bmatrix}$$

5.1.3 Deriving $\frac{\partial \vec{z}^{[k]T}}{\partial w_j^{[k]}}$ and $\frac{\partial \vec{z}^{[k]T}}{\partial b^{[k]}}$

$$\frac{\partial \vec{z}^{[k]T}}{\partial w_j^{[k]}} = \begin{bmatrix} \frac{\partial z^{[k](1)}}{\partial w_j^{[k]}} \\ \frac{\partial z^{[k](2)}}{\partial w_j^{[k]}} \\ \vdots \\ \frac{\partial z^{[k](m)}}{\partial w_j^{[k]}} \end{bmatrix}, \quad \frac{\partial \vec{z}^{[k]T}}{\partial b^{[k]}} = \begin{bmatrix} \frac{\partial z^{[k](1)}}{\partial b^{[k]}} \\ \frac{\partial z^{[k](2)}}{\partial b^{[k]}} \\ \vdots \\ \frac{\partial z^{[k](m)}}{\partial b^{[k]}} \end{bmatrix}$$

For $i = 1, 2, \dots, m$, and $j = 1, 2, \dots, n_x$, we have the following expressions:

$$\frac{\partial z^{[k](i)}}{\partial w_j^{[k]}} = \frac{\partial (\bar{w}^{[k]T} \vec{x}^{(i)} + b)}{\partial w_j^{[k]}} = \frac{\partial \left(\left(\sum_{j=1}^{n_x} w_j^{[k]} x_j^{(i)} \right) + b \right)}{\partial w_j^{[k]}} = x_j^{(i)}$$

$$\frac{\partial z^{[k](i)}}{\partial b^{[k]}} = \frac{\partial (\bar{w}^{[k]T} \vec{x}^{(i)} + b)}{\partial b^{[k]}} = \frac{\partial \left(\left(\sum_{j=1}^{n_x} w_j^{[k]} x_j^{(i)} \right) + b \right)}{\partial b^{[k]}} = 1$$

Thus, we have the following results:

$$\frac{\partial \vec{z}^{[k]T}}{\partial w_j^{[k]}} = \begin{bmatrix} x_j^{(1)} \\ x_j^{(2)} \\ \vdots \\ x_j^{(m)} \end{bmatrix}, \quad \frac{\partial \vec{z}^{[k]T}}{\partial b^{[k]}} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \mathbf{1}_{m \times 1}$$

5.1.4 Putting It All Together

Using the derivations above, we have the following results:

5.1.4.1 $\frac{\partial \mathcal{J}^{[k]}}{\partial w_j^{[k]}}$

$$\begin{aligned} \frac{\partial \mathcal{J}^{[k]}}{\partial w_j^{[k]}} &= \frac{\partial \mathcal{J}^{[k]}}{\partial \vec{a}^{[k]}} \times \frac{\partial \vec{a}^{[k]}}{\partial \vec{z}^{[k]T}} \times \frac{\partial \vec{z}^{[k]T}}{\partial w_j^{[k]}} \\ &= \frac{1}{m} \left[\frac{a^{(1)[k]} - y^{(1)}}{a^{(1)[k]} \times (1 - a^{(1)[k]})}, \frac{a^{(2)[k]} - y^{(2)}}{a^{(2)[k]} \times (1 - a^{(2)[k]})}, \dots, \frac{a^{(m)[k]} - y^{(m)}}{a^{(m)[k]} \times (1 - a^{(m)[k]})} \right] \\ &\quad \times \begin{bmatrix} a^{(1)[k]} \times (1 - a^{(1)[k]}) & 0 & \dots & 0 \\ 0 & a^{(2)[k]} \times (1 - a^{(2)[k]}) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & a^{(m)[k]} \times (1 - a^{(m)[k]}) \end{bmatrix} \\ &\quad \times \begin{bmatrix} x_j^{(1)} \\ x_j^{(2)} \\ \vdots \\ x_j^{(m)} \end{bmatrix} \\ &= \frac{1}{m} [a^{(1)[k]} - y^{(1)}, a^{(2)[k]} - y^{(2)}, \dots, a^{(m)[k]} - y^{(m)}] \times \begin{bmatrix} x_j^{(1)} \\ x_j^{(2)} \\ \vdots \\ x_j^{(m)} \end{bmatrix} \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{m} ([a^{(1)[k]}, a^{(2)[k]}, \dots, a^{(m)[k]}] - [y^{(1)}, y^{(2)}, \dots, y^{(m)}]) \times \begin{bmatrix} x_j^{(1)} \\ x_j^{(2)} \\ \vdots \\ x_j^{(m)} \end{bmatrix} \\
&= \frac{1}{m} (A^{[k]} - Y) \times \begin{bmatrix} x_j^{(1)} \\ x_j^{(2)} \\ \vdots \\ x_j^{(m)} \end{bmatrix} \\
&= \frac{1}{m} \times [x_j^{(1)}, x_j^{(2)}, \dots, x_j^{(m)}] \times (A^{[k]} - Y)^T
\end{aligned}$$

In vector form, we have the following result:

$$\begin{aligned}
\frac{\partial \mathcal{J}^{[k]}}{\partial \vec{w}^{[k]T}} &= \begin{bmatrix} \frac{\partial \mathcal{J}^{[k]}}{\partial w_1^{[k]}} \\ \frac{\partial \mathcal{J}^{[k]}}{\partial w_2^{[k]}} \\ \vdots \\ \frac{\partial \mathcal{J}^{[k]}}{\partial w_{n_x}^{[k]}} \end{bmatrix} = \begin{bmatrix} \frac{1}{m} \times [x_1^{(1)}, x_1^{(2)}, \dots, x_1^{(m)}] \times (A^{[k]} - Y)^T \\ \frac{1}{m} \times [x_2^{(1)}, x_2^{(2)}, \dots, x_2^{(m)}] \times (A^{[k]} - Y)^T \\ \vdots \\ \frac{1}{m} \times [x_{n_x}^{(1)}, x_{n_x}^{(2)}, \dots, x_{n_x}^{(m)}] \times (A^{[k]} - Y)^T \end{bmatrix} \\
&= \frac{1}{m} \times \begin{bmatrix} [x_1^{(1)}, x_1^{(2)}, \dots, x_1^{(m)}] \\ [x_2^{(1)}, x_2^{(2)}, \dots, x_2^{(m)}] \\ \vdots \\ [x_{n_x}^{(1)}, x_{n_x}^{(2)}, \dots, x_{n_x}^{(m)}] \end{bmatrix} \times (A^{[k]} - Y)^T \\
&= \frac{1}{m} \times \begin{bmatrix} [x_1^{(1)}] \\ [x_2^{(1)}] \\ \vdots \\ [x_{n_x}^{(1)}] \end{bmatrix} \begin{bmatrix} [x_1^{(2)}] \\ [x_2^{(2)}] \\ \vdots \\ [x_{n_x}^{(2)}] \end{bmatrix} \dots \begin{bmatrix} [x_1^{(m)}] \\ [x_2^{(m)}] \\ \vdots \\ [x_{n_x}^{(m)}] \end{bmatrix} \times (A^{[k]} - Y)^T \\
&= \frac{1}{m} \times [\vec{x}^{(1)} \quad \vec{x}^{(2)} \quad \dots \quad \vec{x}^{(m)}] \times (A^{[k]} - Y)^T \\
&= \frac{1}{m} \times X \times (A^{[k]} - Y)^T
\end{aligned}$$

5.1.4.2 $\frac{\partial \mathcal{J}^{[k]}}{\partial b^{[k]}}$

Leveraging previous derivations, we have the following result:

$$\begin{aligned}
\frac{\partial \mathcal{J}^{[k]}}{\partial b^{[k]}} &= \frac{\partial \mathcal{J}^{[k]}}{\partial \vec{a}^{[k]}} \times \frac{\partial \vec{a}^{[k]}}{\partial \vec{z}^{[k]T}} \times \frac{\partial \vec{z}^{[k]T}}{\partial b^{[k]}} = \frac{1}{m} [a^{(1)[k]} - y^{(1)}, a^{(2)[k]} - y^{(2)}, \dots, a^{(m)[k]} - y^{(m)}] \times \mathbf{1}_{m \times 1} \\
&= \frac{1}{m} \sum_{i=1}^m a^{(i)[k]} - y^{(i)}
\end{aligned}$$

5.2 Checking Direction of Gradient Descent

We conduct a second derivative test on each coefficient:

$$\begin{aligned}
\frac{\partial \mathcal{J}^2[k]}{\partial w_j^{[k]^2}} &= \frac{\partial \left(\frac{\partial \mathcal{J}^{[k]}}{\partial w_j^{[k]}} \right)}{\partial w_j^{[k]}} = \frac{\partial \left(\frac{1}{m} [a^{(1)[k]} - y^{(1)}, a^{(2)[k]} - y^{(2)}, \dots, a^{(m)[k]} - y^{(m)}] \times \begin{bmatrix} x_j^{(1)} \\ x_j^{(2)} \\ \vdots \\ x_j^{(m)} \end{bmatrix} \right)}{\partial w_j^{[k]}} \\
&= \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \times \frac{\partial a^{(i)[k]}}{\partial w_j^{[k]}} \\
&= \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \times \frac{\partial a^{(i)[k]}}{\partial z^{(i)[k]}} \times \frac{\partial z^{(i)[k]}}{\partial w_j^{[k]}} \\
&= \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \times (a^{(i)[k]} \times (1 - a^{(i)[k]})) \times x_j^{(i)} \\
&= \frac{1}{m} \sum_{i=1}^m (x_j^{(i)})^2 \times (a^{(i)[k]} \times (1 - a^{(i)[k]}))
\end{aligned}$$

Note that since $a^{(i)[k]} = g(z^{(i)[k]}) = \frac{1}{1+e^{-z^{(i)[k]}}}$, therefore $a^{(i)[k]} \in (0,1)$; thus $1 - a^{(i)[k]} \in (0,1)$. Finally, since $m > 0$ and $(x_j^{(i)})^2 \geq 0$, if we assume that there is at least one $x_j^{(i)} > 0$, we thus have $\frac{\partial \mathcal{J}^2[k]}{\partial w_j^{[k]^2}} > 0$ ensuring that $\mathcal{J}^{[k]}$ is strictly convex with respect $w_j^{[k]}$ for $j = 1, 2, \dots, n_x$.

$$\begin{aligned}
\frac{\partial \mathcal{J}^2[k]}{\partial b^{[k]^2}} &= \frac{\partial \left(\frac{\partial \mathcal{J}^{[k]}}{\partial w_j^{[k]}} \right)}{\partial b^{[k]}} = \frac{\partial \left(\frac{1}{m} \sum_{i=1}^m a^{(i)[k]} - y^{(i)} \right)}{\partial b^{[k]}} = \frac{1}{m} \sum_{i=1}^m \frac{\partial a^{(i)[k]}}{\partial b^{[k]}} \\
&= \frac{1}{m} \sum_{i=1}^m \frac{\partial a^{(i)[k]}}{\partial z^{(i)[k]}} \times \frac{\partial z^{(i)[k]}}{\partial b^{[k]}} \\
&= \frac{1}{m} \sum_{i=1}^m (a^{(i)[k]} \times (1 - a^{(i)[k]})) \times 1 > 0
\end{aligned}$$

Thus, $\frac{\partial \mathcal{J}^2[k]}{\partial b^{[k]^2}} > 0$, ensuring that $\mathcal{J}^{[k]}$ is strictly convex with respect $b^{[k]}$ regardless of the values of $x_j^{(i)}$.

5.3 Final Derivations

- \vec{w} (more efficient to express in vector form):
 - $\vec{w}^{[0]} = \text{vector of random numbers (typically chosen to be close to zero)}$
 - $\vec{w}^{[k+1]} = \vec{w}^{[k]} - \alpha \times \left(\frac{1}{m} \times X \times (A^{[k]} - Y)^T \right)$, where α is a set learning rate (e.g., 0.005), and $k = 0, 1, \dots, N$, where N is the total number of iterations
- b :
 - $b^{[0]} = \text{random number (typically chosen to be zero)}$
 - $b^{[k+1]} = b^{[k]} - \alpha \times \left(\frac{1}{m} \sum_{i=1}^m a^{(i)[k]} - y^{(i)} \right)$, where α is a set learning rate (e.g., 0.005), and $k = 0, 1, \dots, N$, where N is the total number of iterations