

# Logistic Regression: back-propagation derivation

Harsha Vardhan

February 6, 2022

## Abstract

This document contains derivation of the gradients for a logistic-regression classifier, using back-propagation. For the implementation of the classifier, see the accompanying notebooks.

## 1 Network Architecture

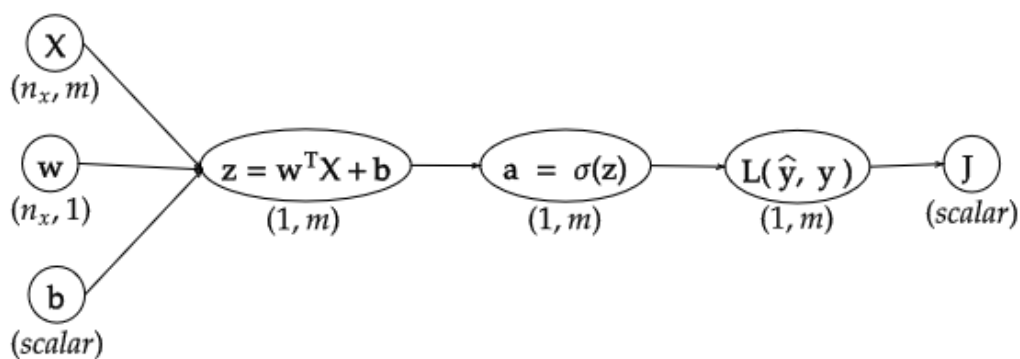


Figure 1: the computation graph along with the dimensions of each nodes' output

## 2 Forward Propagation

The equations for forward propagation are as follows:

$$\begin{aligned}\mathbf{z} &= \mathbf{w}^T \mathbf{X} + b \vec{1}_{(1,m)} \\ \mathbf{a} = \hat{\mathbf{y}} &= \frac{1}{1 + e^{-\mathbf{z}}} \\ \mathbf{L}(\hat{\mathbf{y}}, \mathbf{y}) &= -\mathbf{y} \log(\hat{\mathbf{y}}) - (1 - \mathbf{y}) \log(1 - \hat{\mathbf{y}}) \\ J &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{m} \mathbf{L} \vec{1}_{(m,1)}\end{aligned}$$

where,  $m$  is the number of training-examples  
 $\mathbf{w}$  is a  $(n_x, 1)$  dimensional vector  
 $b$  is a scalar-value  
 $\mathbf{L}$  is the loss function, and is a  $(1, m)$  dimensional vector  
 $J$  is the cost-function

## 3 Optimization: gradient-descent

The optimization is performed according to the following equations:

$$\begin{aligned}\mathbf{w} &:= \mathbf{w} - \alpha \nabla_{\mathbf{w}}(J) \\ b &:= b - \alpha \nabla_b(J)\end{aligned}$$

where,  $\alpha$  is the learning-rate/step-size.

### 3.1 Back-propagation

The gradients in the above equations are derived using back-propagation, as follows:

Since,  $J$  is a scalar, we can write  $J = \text{tr}(J) = J^\top = \text{tr}(J^\top)$ . And the derivative can be computed as follows:

$$\begin{aligned} dJ &= d(\text{tr}(J)) \\ &= \text{tr}(dJ) \end{aligned}$$

The objective while computing the above derivative, is to massage the expression to the following form:

$$dy = \text{tr}(\mathbf{A}d\mathbf{X})$$

then,

$$\frac{dy}{d\mathbf{X}} = \mathbf{A}$$

See [1] and [2] for more information.

### 3.1.1 Computing $\frac{dJ}{d\mathbf{L}^\top}$

$$\begin{aligned} dJ &= \text{tr}(dJ^\top) \\ &= \text{tr}\left(d\left(\frac{1}{m}(\vec{\mathbf{1}}_{(m,1)})^\top \mathbf{L}^\top\right)\right) \\ &= \text{tr}\left(\frac{1}{m}(\vec{\mathbf{1}}_{(m,1)})^\top d\mathbf{L}^\top\right) \\ \implies \frac{dJ}{d\mathbf{L}^\top} &= \frac{1}{m}(\vec{\mathbf{1}}_{(m,1)})^\top \end{aligned}$$

### 3.1.2 Computing $\frac{dJ}{d\mathbf{a}^\top}$

$$\begin{aligned} dJ &= \text{tr}\left(\frac{dJ}{d\mathbf{L}^\top} \frac{d\mathbf{L}^\top}{d\mathbf{a}^\top} d\mathbf{a}^\top\right) \\ &= \text{tr}\left(\frac{1}{m}(\vec{\mathbf{1}}_{(m,1)})^\top \begin{bmatrix} \frac{\partial L^{(1)}}{\partial a^{(1)}} & 0 & \cdots & 0 \\ 0 & \frac{\partial L^{(2)}}{\partial a^{(2)}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{\partial L^{(m)}}{\partial a^{(m)}} \end{bmatrix} d\mathbf{a}^\top\right) \\ \implies \frac{dJ}{d\mathbf{a}^\top} &= \frac{1}{m} \left[ \frac{\partial L^{(1)}}{\partial a^{(1)}}, \frac{\partial L^{(2)}}{\partial a^{(2)}}, \dots, \frac{\partial L^{(m)}}{\partial a^{(m)}} \right] \end{aligned}$$

where,

$$\begin{aligned}\frac{d\mathbf{L}^\top}{d\mathbf{a}^\top} &= \left[ \frac{\partial}{\partial a^{(1)}}, \frac{\partial}{\partial a^{(2)}}, \dots, \frac{\partial}{\partial a^{(m)}} \right] \otimes \begin{bmatrix} L^{(1)} \\ L^{(2)} \\ \vdots \\ L^{(m)} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial L^{(1)}}{\partial a^{(1)}} & 0 & \dots & 0 \\ 0 & \frac{\partial L^{(2)}}{\partial a^{(2)}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{\partial L^{(m)}}{\partial a^{(m)}} \end{bmatrix}\end{aligned}$$

### 3.1.3 Computing $\frac{dJ}{dz^\top}$

$$\begin{aligned}dJ &= \text{tr} \left( \frac{dJ}{d\mathbf{a}^\top} \frac{d\mathbf{a}^\top}{dz^\top} dz^\top \right) \\ &= \text{tr} \left( \frac{dJ}{d\mathbf{a}^\top} \begin{bmatrix} \frac{\partial a^{(1)}}{\partial z^{(1)}} & 0 & \dots & 0 \\ 0 & \frac{\partial a^{(2)}}{\partial z^{(2)}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{\partial a^{(m)}}{\partial z^{(m)}} \end{bmatrix} dz^\top \right) \\ &= \text{tr} \left( \frac{dJ}{d\mathbf{a}^\top} \circ \text{diag}^{-1} \left( \frac{d\mathbf{a}^\top}{dz^\top} \right) dz^\top \right) \\ \Rightarrow \frac{dJ}{dz^\top} &= \frac{dJ}{d\mathbf{a}^\top} \circ \text{diag}^{-1} \left( \frac{d\mathbf{a}^\top}{dz^\top} \right)\end{aligned}$$

where,

$$\begin{aligned}\frac{d\mathbf{a}^\top}{dz^\top} &= \left[ \frac{\partial}{\partial z^{(1)}}, \frac{\partial}{\partial z^{(2)}}, \dots, \frac{\partial}{\partial z^{(m)}} \right] \otimes \begin{bmatrix} a^{(1)} \\ a^{(2)} \\ \vdots \\ a^{(m)} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial a^{(1)}}{\partial z^{(1)}} & 0 & \dots & 0 \\ 0 & \frac{\partial a^{(2)}}{\partial z^{(2)}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{\partial a^{(m)}}{\partial z^{(m)}} \end{bmatrix}\end{aligned}$$

### 3.1.4 Computing $\frac{dJ}{d\mathbf{w}}$ , and $\frac{dJ}{db}$

$$dJ = \text{tr}\left(\frac{dJ}{d\mathbf{z}^\top} d\mathbf{z}^\top\right)$$

where,  $d\mathbf{z}^\top$  may be expanded as,

$$d\mathbf{z}^\top = d\mathbf{X}^\top \mathbf{w} + \mathbf{X}^\top d\mathbf{w} + (\vec{1}_{(1,m)})^\top db$$

Here,  $d\mathbf{X}^\top = 0$ , since,  $\mathbf{X}$  is the input. And when differentiating w.r.t  $\mathbf{w}$ , we have  $db = 0$ . So,

$$\begin{aligned} dJ &= \text{tr}\left(\frac{dJ}{d\mathbf{z}^\top} \mathbf{X}^\top d\mathbf{w}\right) \\ \implies \frac{dJ}{d\mathbf{w}} &= \frac{dJ}{d\mathbf{z}^\top} \mathbf{X}^\top \end{aligned}$$

when differentiating w.r.t  $b$ , we have  $d\mathbf{w} = 0$ . So,

$$\begin{aligned} dJ &= \text{tr}\left(\frac{dJ}{d\mathbf{z}^\top} (\vec{1}_{(1,m)})^\top db\right) \\ \implies \frac{dJ}{db} &= \frac{dJ}{d\mathbf{z}^\top} (\vec{1}_{(1,m)})^\top \end{aligned}$$

## 3.2 Jacobian or Gradient?

In the above derivations, we have used the numerator layout while performing matrix-derivatives. Therefore, the derivatives are jacobians and not gradients. And the corresponding gradients are simply transpose of jacobians, i.e.,

$$\begin{aligned} \nabla_{\mathbf{w}}(J) &= \left(\frac{dJ}{d\mathbf{w}}\right)^\top = \mathbf{X} \left(\frac{dJ}{d\mathbf{z}^\top}\right)^\top \\ \nabla_b(J) &= \left(\frac{dJ}{db}\right)^\top = \vec{1}_{(1,m)} \left(\frac{dJ}{d\mathbf{z}^\top}\right)^\top \end{aligned}$$

## References

- [1] T. Minka, "Old and new matrix algebra useful for statistics," Sep. 1997. [Online]. Available: <https://www.microsoft.com/en-us/research/publication/old-new-matrix-algebra-useful-statistics/>.
- [2] J. R. Magnus and H. Neudecker, *Matrix Differential Calculus with Applications in Statistics and Econometrics*, Second. John Wiley, 1999, ISBN: 0471986321 9780471986324 047198633X 9780471986331.