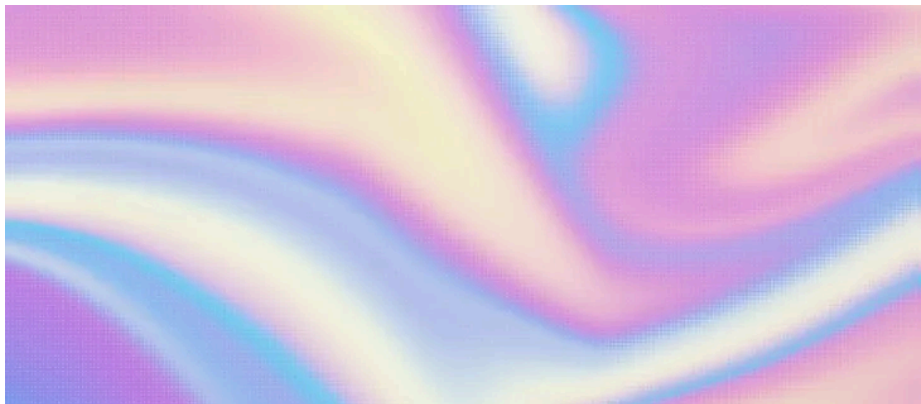G Gianluca Turcatel ♛    Dec 23, 2021    1 min read

# Derivation of the Binary Cross Entropy Loss Gradient



The binary cross entropy loss function is the preferred loss function in binary classification tasks, and is utilized to estimate the value of the model's parameters through gradient descent. In order to apply gradient descent we must calculate the derivative (gradient) of the loss function w.r.t. the model's parameters. Deriving the gradient is usually the most tedious part of training a machine learning model.
In this article we will derive the derivative of the binary cross entropy loss function w.r.t. W, step by step.

The  binary cross entropy loss is given by

$$L(W) = -[y * \log(\hat{y}) + (1 - y) * \log(1 - \hat{y})] \quad (1)$$

y is the observed class, y_hat the prediction, W the model's parameters. Predictions are given by:

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

z is equal to:

$$z = w_0 + x_1 * w_1 + x_2 * w_2 + \cdots + x_k * w_k = W^T X \quad (3)$$

To calculate the gradient of L(W) w.r.t. W we will use the chain rule:

$$\frac{\partial L(W)}{\partial W} = \frac{\partial L(W)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial W} \quad (4)$$

Let's derive the first term:

$$\frac{\partial L(W)}{\partial \hat{y}} = -\left(\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}}\right) = \frac{\hat{y} - y}{\hat{y}(1-\hat{y})} \qquad (5)$$

The second term is a little more complicated:

$$\frac{\partial \hat{y}}{\partial z} = \frac{\partial}{\partial z}\left[\frac{1}{1+e^{-z}}\right] = \frac{\partial}{\partial z}(1+e^{-z})^{-1} = \qquad (6)$$

$$= -(1+e^{-z})^{-2}(-e^{-z}) = \frac{e^{-z}}{(1+e^{-z})^2} = \qquad (7)$$

$$= \frac{1}{(1+e^{-z})}\frac{e^{-z}}{(1+e^{-z})} = \frac{1}{1+e^{-z}}\frac{(1+e^{-z})-1}{(1+e^{-z})} = (8)$$

$$= \frac{1}{1 + e^{-z}} \left( \frac{1 + e^{-z}}{1 + e^{-z}} - \frac{1}{1 + e^{-z}} \right) = \qquad (9)$$

$$= \frac{1}{1 + e^{-z}} \left( 1 - \frac{1}{1 + e^{-z}} \right) = \hat{y}(1 - \hat{y}) \qquad (10)$$

Done with the second term. The derivative of the third term is straight forward:

$$\frac{\partial z}{\partial W} = X \qquad (11)$$

Now let's put everything together:

$$\frac{\partial L(W)}{\partial W} = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \hat{y}(1 - \hat{y})X = (\hat{y} - y)X \qquad (12)$$

And there you have it: the derivative of the binary cross entropy loss function w.r.t. the model's parameters.

Follow me on Twitter and Facebook to stay updated.