

Binary classification w/ logistic regression

$$z = w^T x + b$$

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

will be between 0 and 1

logistic function.

$y > 1/2 \rightarrow$  positive  
 $y < 1/2 \rightarrow$  negative.

target  $\in \{0, 1\}$

Loss (cross-entropy):  $\mathcal{L}_{CE} = -t \log y - (1-t) \log(1-y)$   
model's prediction

Problem: when  $z \ll 0$  (model is very confident that  $x$  is a negative example).

$$z \rightarrow -\infty$$

$$e^{-z} \rightarrow +\infty$$

$$y = \sigma(z) = \frac{1}{1 + e^{-z}} \rightarrow \frac{1}{\infty} \rightarrow 0$$

logistic function

$-\log y = -\log 0$  tries to compute  $\log(0)$   
 $-\infty$

PB: Hard-to-find subtle bugs because we're trying to compute  $\log(0)$ .

Solution: combine the loss function (cross-entropy) with the activation function (logistic function).

Logistic-cross-entropy function:

$$\mathcal{L}_{LCE}(z, t) = \mathcal{L}_{CE}(\sigma(z), t) = -t \log(\sigma(z)) - (1-t) \log(1-\sigma(z))$$

$z = w^T x + b$  rather than model's output  
 $y = \sigma(z)$

target  $t \in \{0, 1\}$

$$\begin{aligned} &= -t \log\left(\frac{1}{1+e^{-z}}\right) - (1-t) \log\left(1 - \frac{1}{1+e^{-z}}\right) \\ &= -t \left( \log(1) - \log(1+e^{-z}) \right) - (1-t) \log\left(\frac{1+e^{-z}-1}{1+e^{-z}}\right) \\ &= +t \log(1+e^{-z}) - (1-t) \log\left(\frac{e^{-z}}{1+e^{-z}}\right) \\ &= t \log(1+e^{-z}) - (1-t) \log\left(\frac{1}{e^z+1}\right) \\ &= t \log(1+e^{-z}) - (1-t) \left( \log(1) - \log(e^z+1) \right) \end{aligned}$$

cross-entropy

Combined loss and activation function

$$\mathcal{L}_{LCE}(z, t) = t \log(1+e^{-z}) + (1-t) \log(1+e^z)$$

logistic

$\neq 0$   $\neq 0$

More numerically stable.

NumPy Implementation

$$E = t * \text{np.logaddexp}(0, -z) + (1-t) * \text{np.logaddexp}(0, z)$$

Combined loss and activation function

function included in NumPy which computes

$$\text{np.logaddexp}(a, b) = \log(e^a + e^b)$$

np here stands for NumPy  
import numpy as np

Gradient descent

Update rule:  $w \leftarrow w - \alpha \frac{\partial \mathcal{L}}{\partial w}$

learning rate (hyperparameter)

For simplicity, assume that  $w$  is a scalar but using vectorized notation, we could derive a similar result.

$\mathcal{L}_{LCE}$  combines loss and activation functions

$$\frac{\partial \mathcal{L}(z, t)}{\partial w} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial w}$$

using chain rule.

per definition of  $z$

$$\frac{\partial z}{\partial w} = \frac{\partial (wx + b)}{\partial w} = x$$

$$\begin{aligned} \frac{\partial \mathcal{L}_{LCE}(z, t)}{\partial z} &= \frac{\partial (t \log(1+e^{-z}) + (1-t) \log(1+e^z))}{\partial z} \\ &= \frac{\partial (t \log(1+e^{-z}))}{\partial z} + \frac{\partial ((1-t) \log(1+e^z))}{\partial z} \\ &= t \frac{\partial \log(1+e^{-z})}{\partial z} + (1-t) \frac{\partial \log(1+e^z)}{\partial z} \\ &= t \frac{(-1)e^{-z}}{1+e^{-z}} + (1-t) \frac{e^z}{1+e^z} \times e^z \\ &= -t \left( \frac{1+e^{-z}}{1+e^{-z}} - \frac{1}{1+e^{-z}} \right) + (1-t) \frac{1}{e^{-z}+1} \\ &= -t \left( \frac{1+e^{-z}}{1+e^{-z}} - \frac{1}{1+e^{-z}} \right) + (1-t) y \\ &= -t \left( 1 - \frac{1}{1+e^{-z}} \right) + (1-t) y \\ &= -t(1-y) + (1-t)y \\ &= -t + y + y - ty \\ &= y - t \end{aligned}$$

logistic function  $y = \sigma(z) = \frac{1}{1+e^{-z}}$

$$\frac{\partial \mathcal{L}_{LCE}(z, t)}{\partial z} = y - t$$

$$\frac{\partial \mathcal{L}_{LCE}(z, t)}{\partial w} = \frac{\partial \mathcal{L}_{LCE}(z, t)}{\partial z} \frac{\partial z}{\partial w} = (y-t)x$$

Update rule for binary linear classifier with a logistic activation function and a cross-entropy loss function

$$w \leftarrow w - \alpha (y-t)x$$

If model predicts positive but label is negative ( $t=0$ )  $y > t$

$$y > t \Rightarrow y - t > 0$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial z} > 0$$

is consistent w/ the intuition that to decrease the loss (improve the prediction) we need to decrease  $z$



If model predicts negative but label is positive ( $t=1$ )  $y < t$

$$y < t \Rightarrow y - t < 0$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial z} < 0$$

corresponds to our intuition that we should increase  $z$  to decrease the loss

