

$$\hat{y} = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

Given $(x_1, y_1) \dots$

Given m datapoints $(\langle x_1^1, \dots, x_n^1 \rangle, y^1) \dots$

$\dots (\langle x_1^m, \dots, x_n^m \rangle, y^m)$

classifier

→ This does not work for logistic regression

$$\hat{y} = \frac{1}{1 + e^{-(w_1 x_1 + \dots + w_n x_n + b)}}$$

$$= \sigma(w_1 x_1 + \dots + w_n x_n + b)$$

Sigmoid function

Cost function

$$\text{cost} = \frac{1}{m} \sum_{i=0}^m |y - \hat{y}|$$

This is suitable for Linear regression

In classification $y = 0$ or 1

↳ If we take the above formula, we may end up in local minima instead of global minima.

(Not proving this. pls take it as a fact)



So, we use Binary Cross Entropy

$$\text{cost} = -\frac{1}{m} [y \log \hat{y} + (1-y) \log (1-\hat{y})]$$

$$\hat{y} = \sigma(W^T x + b)$$

How is the above equation a cost function?
Let's try to answer that.

When $y \neq \hat{y}$ are not the same, cost should be high. Let's check that.

$$\text{cost} = -\frac{1}{m} \sum [y \log \hat{y} + (1-y) \log (1-\hat{y})]$$

$\epsilon = \text{epsilon}$

$y = 0$	\Rightarrow	0	$-1 \log (1-\epsilon)$	\uparrow close to 1 very small quantity
$\hat{y} \approx 0$				
$y = 0$		0	$-1 \log (0+\epsilon)$	\uparrow close to zero very high quantity
$\hat{y} \approx 1$				
$y = 1$		$-\log (0+\epsilon)$	0	
$\hat{y} \approx 0$		close to 0 very high quantity		
$y = 1$		$-\log (1-\epsilon)$	0	
$\hat{y} \approx 1$		close to 1 very small quantity		

Conclusion: The formula works well as cost function

Derivative of the cost.

$$\text{cost} = -\frac{1}{m} \sum_{i=1}^m y \log \hat{y} + (1-y) \log (1-\hat{y})$$

$$\hat{y} = \sigma(W^T X + b)$$

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

$$z = W^T X + b$$

$$W = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

To find $\frac{\partial \text{cost}}{\partial w} + \frac{\partial \text{cost}}{\partial b}$

$$\begin{bmatrix} \frac{\partial \text{cost}}{\partial w_1} \\ \vdots \\ \frac{\partial \text{cost}}{\partial w_n} \end{bmatrix}$$

$$\begin{aligned}\frac{\partial \text{cost}}{\partial w_1} &= \frac{\partial \text{cost}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_1} \\ &= \frac{\partial \text{cost}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_1}\end{aligned}$$

$$\text{cost} = -\frac{1}{m} \sum_{i=1}^m y \log \hat{y} + (1-y) \log (1-\hat{y})$$

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

$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

$$\begin{aligned}\frac{\partial \text{cost}}{\partial \hat{y}} &= -\frac{y}{\hat{y}} + \frac{(1-y)}{(1-\hat{y})} \\ &= \frac{-y + y\hat{y} + \hat{y} - y\hat{y}}{\hat{y}(1-\hat{y})} \\ &= \frac{-y + \hat{y}}{\hat{y}(1-\hat{y})}\end{aligned}$$

$$\begin{aligned}\frac{\partial \hat{y}}{\partial z} &= \frac{e^{-z}}{(1+e^{-z})^2} \\ &= \frac{(1-\hat{y})}{\hat{y}} \hat{y}^2 \\ &= (1-\hat{y}) \hat{y}\end{aligned}$$

$$\frac{\partial z}{\partial w_1} = x_1$$

$$\begin{aligned}\frac{\partial \text{cost}}{\partial w_1} &= \frac{-y + \hat{y}}{\hat{y}(1-\hat{y})} \cdot (1-\hat{y}) \hat{y} \cdot x_1 \\ &= (\hat{y} - y) x_1\end{aligned}$$

$$\begin{aligned}\frac{\partial x^n}{\partial x} &= n x^{n-1} \\ \frac{\partial}{\partial x} \left(\frac{1}{x} \right) &= -\frac{1}{x^2} \\ \frac{\partial}{\partial x} (\log x) &= \frac{1}{x} \\ \frac{\partial}{\partial x} e^x &= e^x \\ \frac{\partial}{\partial x} [f(x)]^n &= n [f(x)]^{n-1} \times \frac{\partial f(x)}{\partial x}\end{aligned}$$

$$\frac{\partial z}{\partial b} = \frac{\partial \cos t}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial b}$$

$$= (\hat{y} - y) \times 1$$

$$\therefore W = W - 2 \cdot \frac{\partial \text{Cost}}{\partial w}$$

 α = Learning rate.

$$X = [I, I, \dots]$$

Learning-rate = 0.1

Repeat $\binom{N}{1000}$ times

Determining differential of cost w.r.t. w_i 's & b

Implementation in github