

Logistic Regression (Part II)

Till now we have seen logistic regression for binary classification, that uses Sigmoid function

Let's consider the following dataset,

Weight (x_1)	color (x_2)	fruit (y)
150	Red	Apple
10	green	Grape
120	yellow	Banana
8	green	Grape
200	Red	Apple

features

target with 3 classes

Apple, Grape and banana.

We can convert this dataset into (by doing one hot encoding) \rightarrow

2nd sample has class grape, that's why $y_2 = 1$ and $y_1 = y_3 = 0$.

x_1	x_2	$y_1 = \text{Apple}$	$y_2 = \text{Grape}$	$y_3 = \text{Banana}$
150	Red	1	0	0
10	green	0	1	0
120	yellow	0	0	1
8	green	0	1	0
200	Red	1	0	0

features

Class is Apple, so

$y_1 = 1$ and y_2 and y_3 is 0.

As we have three classes here, now we will consider each classes z-score for each classes.

so, we will 3 ~~at~~ z-score equation,

$$z_1 = w_{11}x_1 + w_{12}x_2 + \cancel{w_{13}x_3} + b_1 \rightarrow \text{for } y_1 \text{ column}$$

$$z_2 = w_{21}x_1 + w_{22}x_2 + b_2 \rightarrow \text{for } y_2 \text{ column}$$

$$z_3 = w_{31}x_1 + w_{32}x_2 + b_3 \rightarrow \text{for } y_3 \text{ column.}$$

w_{32} \rightarrow 3rd z score
 \rightarrow weight for second feature x_2 .

so, if we have K classes will have z_1, z_2, \dots, z_K equations. If we have n samples in the dataset, and features.

	x_1, \dots, x_f	y_1	y_2	\dots	y_K
$x^{(1)}$					
$x^{(2)}$					
:					
$x^{(n)}$					

we will have K number of z equations \rightarrow
 $z_1 = w_{11}x_1 + w_{12}x_2 + \dots + w_{1f}x_f + b_1$ \rightarrow f number of features.

$$z_2 = w_{21}x_1 + w_{22}x_2 + \dots + w_{2f}x_f + b_2$$

$$\vdots$$

$$z_K = w_{K1}x_1 + w_{K2}x_2 + \dots + w_{Kf}x_f + b_K$$

Again our goal is to determine these 2 equations for the model. For that we need to find the weights and biases.

Let us we have one sample only, with ~~4~~^{four} features and three classes like below →

	x_1	x_2	x_3	x_4	y_1	\hat{y}_2	\hat{y}_3
sample f	1	0.5	2.0	3.0	?	?	?

features ~~Actual~~ target prediction

As we have three classes, 2 scores →

$$z_1 = w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + w_{14}x_4 + b_1$$

$$z_2 = w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + w_{24}x_4 + b_2$$

$$z_3 = w_{31}x_1 + w_{32}x_2 + w_{33}x_3 + w_{34}x_4 + b_3$$

Let us assume we know all of these weights and bias values. We have to find the prediction for this sample.

Probability of sample being in class $y_1 \rightarrow$

$$\hat{y}_1 = P(y_1 = 1 | x) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

This will give us an prediction probability of the sample being in class 1.

The formula we have used here is Softmax function. This function is used in place of Sigmoid function for multinomial/multi-class Logistic classification.

Similarly, Probability of the sample being class 2/ y_2 =

$$y_2 = P(y_2=1|x) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

and, Probability of the sample being in class 3/ y_3 =

$$y_3 = P(y_3=1|x) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

so, For a vector z of dimensionality K (K classes), Softmax is defined as →

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}, \text{ where } i \in [1, K]$$

we will have →

$$\text{softmax}(z) = \left[\frac{e^{z_1}}{\sum_{j=1}^K e^{z_j}}, \frac{e^{z_2}}{\sum_{j=1}^K e^{z_j}}, \dots, \frac{e^{z_K}}{\sum_{j=1}^K e^{z_j}} \right]$$

If we get,

$$\hat{y}_1 = P(y_1=1|x) = 0.3 \rightarrow \text{predicted class} = 0$$

$$\hat{y}_2 = P(y_2=1|x) = 0.6 \rightarrow \text{predicted class} = 1$$

$$\hat{y}_3 = P(y_3=1|x) = 0.2 \rightarrow \text{predicted class} = 0$$

so, According to this our class will be y_2 as it's probability is largest.

x_1	x_2	x_3	x_4	\hat{y}_1	\hat{y}_2	\hat{y}_3
1	0.5	2.0	3	0	1	0

$$P(y_2=1|x) = 0.6 \xrightarrow{\text{largest}}$$

But while training our model we need to update the weights and biases like we did in logistic regression.

We know to update the parameters we need a loss function. In logistic regression we have used Cross Entropy loss function. We will use it in multiclass classification as well. We know for n samples,

$$L_{CE} = -\sum_{i=1}^n y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})$$

In our previous example → let us assume the Actual classes as below

x_1	x_2	x_3	x_4	y_1	y_2	y_3	\hat{y}_1	\hat{y}_2	\hat{y}_3
1	0.5	2.0	3.0	0	1	0	0.3	0.6	0.2

Actual
values Prediction
Probabilities

The General Loss function in Softmax →

$$L_{CE} = \sum_{i=1}^k y_k \log \hat{y}_k$$

[(1- y_k) term is implicit here because of Softmax function]

In our example $k=3$,

$$\begin{aligned} L_{CE} &= \sum_{i=1}^3 y_k \log \hat{y}_k \\ &= y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + y_3 \log \hat{y}_3 \\ &= 0 \times + 1 \times \log 0.6 + 0 \end{aligned}$$

↓
because y_2 is the correct class here.

The full calculation for LCE hence

$$z_1 = w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + w_{14}x_4 + b_1 = \boxed{w_1x + b_1}$$

$$z_2 = w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + w_{24}x_4 + b_2 = \boxed{w_2x + b_2}$$

$$z_3 = w_{31}x_1 + w_{32}x_2 + w_{33}x_3 + w_{34}x_4 + b_3 = \boxed{w_3x + b_3}$$

$$\hat{y}_2 = P(y_2=1|x) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\therefore L_{CF} = \sum_{k=1}^3 y_k \log \hat{y}_k$$

$$= y_2 \log \hat{y}_2 \rightarrow [y_1 = y_3 = 0]$$

$$= 1 \times \log \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$= \log \frac{e^{w_2x + b_2}}{e^{w_1x + b_1} + e^{w_2x + b_2} + e^{w_3x + b_3}}$$

↓
will have to calculate this term
using all weights and biases.

Now to update the weights using gradient

descent, we will need to use the derivative of the loss function. Let's 3 equations are →

$$z_1 = w_{11}x_1 + w_{12}x_2 + b_1$$

$$z_2 = w_{21}x_1 + w_{22}x_2 + b_2$$

$$z_3 = w_{31}x_1 + w_{32}x_2 + b_3$$

We need to find all of these 9 parameters here.

Let us consider, w_{11} .

We know while updating

$$w_{11}(\text{new}) = w_{11}(\text{old}) - \eta \frac{\delta L_{CE}}{\delta w_{11}}$$

∴ As we have seen in binary Logistic regression

We will use →

$$\frac{\delta L_{CE}}{\delta w_{11}} = (\hat{y}_1 - y_1) x_1$$

Prediction from y_1 class

\hat{y}_1 with x_1

y_1 class (Actual)

b_1 by y_1 class

Similarly we will need all other derivatives,

$$\frac{\delta L_{CE}}{\delta w_{12}}, \frac{\delta L_{CE}}{\delta w_{21}}, \frac{\delta L_{CE}}{\delta w_{22}}, \frac{\delta L_{CE}}{\delta w_{31}}, \frac{\delta L_{CE}}{\delta w_{32}}$$

$$\text{and } \frac{\delta L_{CE}}{\delta b_1}, \frac{\delta L_{CE}}{\delta b_2}, \frac{\delta L_{CE}}{\delta b_3}$$

here $\frac{\delta L_{CE}}{\delta w_{31}} = (\hat{y}_3 - y_3) x_1$

\hat{y}_3 for prediction for y_3
 x_1 weight with w_{31}

And $\frac{\delta L_{CE}}{\delta b_3} = (\hat{y}_3 - y_3)$ bias for \hat{y}_3
 condition. \hat{y}_3 is for y_3 class

Using these derivatives we can update the parameter in Gradient Descent approach.

The general formula for derivative of LCE

with respect to w_{ki}
 \downarrow feature no.
 class no.

$$\frac{\delta L_{CE}}{\delta w_{ki}} = (\hat{y}_k - y_k) x_i$$

and with respect to b_k

$$\frac{\delta L_{CE}}{\delta b_k} = (\hat{y}_k - y_k)$$

Example :

Do one cycle of gradient Descent update for the following dataset.

x_1	x_2	y_1	y_2	y_3
1	2	1	0	0

$$\text{where, } w_{11} = 0.2, w_{12} = (-0.10), w_{21} = 0.0, w_{22} = 0.10$$

$$w_{31} = -0.2, w_{32} = 0.05$$

$$\eta = 0.1, b_1 = 0.10, b_2 = 0.00, b_3 = -0.05$$

Solution: here,

$$z_1 = w_{11}x_1 + w_{12}x_2 + b_1$$

$$z_2 = w_{21}x_1 + w_{22}x_2 + b_2$$

$$z_3 = w_{31}x_1 + w_{32}x_2 + b_3$$

$$\therefore z_1 = (0.2) \times 1 + (-0.10) \times 2 + (0.10) = 0.10$$

$$z_2 = (0.0) \times 1 + (0.10) \times 2 + (0.00) = 0.20$$

$$z_3 = (-0.2) \times 1 + (0.05) \times 2 + (-0.05) = -0.15$$

$$\therefore e^{z_1} = 1.105$$

$$e^{z_2} = 1.221$$

$$e^{z_3} = 0.861$$

$$\text{And } e^{z_1} + e^{z_2} + e^{z_3} = 3.187$$

$$\hat{y}_1 = P(y_1=1) = \text{softmax}(z_1)$$

$$= \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}} = \frac{1 \cdot 105}{3 \cdot 187} \\ = 0.3467$$

$$\hat{y}_2 = P(y_2=1) = \text{softmax}(z_2)$$

$$= \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}} = \frac{1.221}{3.187} \\ = 0.383$$

$$\hat{y}_3 = P(y_3=1) = \text{softmax}(z_3)$$

$$= \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}} = \frac{1.221}{3.187} \frac{0.8607}{3.187} \\ = 0.27$$

As $\hat{y}_2 = 0.383$ is highest, predicted class is

y_2 .

$$LCE = \sum_{i=1}^3 y_i \log \hat{y}_i$$

$$= y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + y_3 \log \hat{y}_3$$

$$= 1 \times \log 0.3467 + 0 \times \log \hat{y}_2 + 0 \log \hat{y}_3$$

$$= 1.059$$

Updates \rightarrow

$$w_{11} = w_{11} - \eta (\hat{y}_1 - y_1) u_1 = 0.2 - 0.1 \times (0.3467 - 1) 1 \\ \approx 0.2653$$

$$w_{12} = w_{12} - \eta (\hat{y}_1 - y_1) u_2 = (-0.10) - 0.1 (0.3467 - 1) 2 \\ = 0.0307$$

$$w_{21} = w_{21} - \eta (\hat{y}_2 - y_2) u_1 = 0.00 - 0.1 (0.3832 - 0) 1 \\ = -0.0383$$

$$w_{22} = w_{22} - \eta (\hat{y}_2 - y_2) u_2 = 0.10 - 0.1 \times (0.3832 - 0) 2 \\ \approx 0.0233$$

$$w_{31} = w_{31} - \eta (\hat{y}_3 - y_3) u_1 = -0.2 - 0.1 \times (0.27 - 0) 1 \\ \approx -0.227$$

$$w_{32} = w_{32} - \eta (\hat{y}_3 - y_3) u_2 = 0.05 - 0.1 (0.27 - 0) 2 \\ = -0.004$$

and, $b_1 = b_1 - \eta (\hat{y}_1 - y_1) = 0.10 - 0.1 \times (0.3467 - 1) \\ \approx 0.1653$

$$b_2 = b_2 - \eta (\hat{y}_2 - y_2) = 0.0 - 0.1 \times (0.3832 - 0) \\ \approx -0.03832$$

$$b_3 = b_3 - \eta (\hat{y}_3 - y_3) = 0.05 - 0.1 \times (0.27 - 0) \\ \approx -0.072$$

Minibatch Gradient Descent (for Binary Classification):

We have already seen an example for stochastic Gradient Descent.

In minibatch process, we divide the dataset into small batches to update the parameters. So, hence updates takes place batch by batch, not sample by sample.

For example →

The diagram illustrates the division of a dataset into minibatches. On the left, a large table represents the original dataset with 5000 samples, having columns for 'features' and 'target'. A red arrow labeled 'divide into 5 batches' points to the right, where five smaller tables are shown, each representing a minibatch. Each minibatch has 1000 samples and is labeled with a circled number (1, 2, 3, 4, or 5) in red. The total number of samples (5000) is also written in red next to the original table.

	features	target
1	:	
:	:	
:	:	
:	:	
5000	:	

5000 samples

	features	target
1000	①	1000
1000	②	1000
1000	③	1000
1000	④	1000
1000	⑤	1000

Each batch has 1000 samples.

The gradient descent will run for each batch and update the parameters after processing each batch. So, the update will take place 5 times here instead of 5000 times like stochastic gradient descent. These 5 batches are called 1 Epoch.

The entire process is ~~repeated~~ repeated for

multiple epochs until the model converges.

Batch Gradient Descent: The entire dataset gets processed at once. So, the parameters gets updated after processing all samples of the dataset.

Example :

Perform one Gradient Descent cycle for the following ~~batch~~ dataset in using minibatch Gradient Descent

x_1	x_2	y
2	1	1
1	3	0
2	2	1

$$w_1 = w_2 = b = 0, \text{ and } \eta = 0.5$$

Solution:

$$z = w_1 x_1 + w_2 x_2 + b$$

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

We know to update the parameters,

$$w_1 = w_1(\text{old}) - \eta \quad \boxed{\frac{\delta LCE}{\delta w_1}} \rightarrow \begin{array}{l} \text{These gradients} \\ \text{will be the} \\ \text{average of individual} \\ \text{gradients.} \end{array}$$
$$w_2 = w_2(\text{old}) - \eta \quad \boxed{\frac{\delta LCE}{\delta w_2}}$$
$$b = b(\text{old}) - \eta \quad \boxed{\frac{\delta LCE}{\delta b}}$$

We know for each sample,

$$\frac{\delta LCE}{\delta w_1} = (\hat{y} - y) x_1 \quad \frac{\delta LCE}{\delta w_2} = (\hat{y} - y) x_2$$

and $\frac{\delta LCE}{\delta b} = (\hat{y} - y)$

But in batch gradient descent it becomes →

$$\frac{\delta LCE}{\delta w_1} = \left(\frac{1}{m} \sum_{i=1}^m [\hat{y}^{(i)} - y^{(i)}] \right) x_1^{(i)}$$

$$\frac{\delta LCE}{\delta w_2} = \left(\frac{1}{m} \sum_{i=1}^m [\hat{y}^{(i)} - y^{(i)}] \right) x_2^{(i)}$$

$$\frac{\delta LCE}{\delta b} = \left(\frac{1}{m} \sum_{i=1}^m [\hat{y}^{(i)} - y^{(i)}] \right)$$

Average of all individual samples' gradients.

cycle - 1 ($w_1 = w_2 = b = 0$)
 $\hat{y} = w_1 x_1 + w_2 x_2 + b$

	x_1	x_2	y	\hat{z}	\hat{y}	$\hat{y} - y$	$(\hat{y} - y)w_1$	$(\hat{y} - y)w_2$
①	2	1	1	0	0.5	-0.5	-1	-0.5
②	1	3	0	0	0.5	0.5	0.5	1.5
③	2	2	1	0	0.5	-0.5	-1	-1

$\hookrightarrow m=3$

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

$$\begin{aligned} \therefore \frac{\delta L}{\delta w_1} &= \frac{1}{m} \sum_{i=1}^m \left[(\hat{y}^{(i)} - y^{(i)}) u^{(i)} \right] \\ &= \frac{1}{m} \left[(\hat{y}^{(1)} - y^{(1)}) u^{(1)} + (\hat{y}^{(2)} - y^{(2)}) u^{(2)} + \right. \\ &\quad \left. (\hat{y}^{(3)} - y^{(3)}) u^{(3)} \right] \\ &= \frac{1}{3} (-1 + 0.5 - 1) = \frac{-1.5}{3} = -0.5 \end{aligned}$$

$$\therefore \frac{\delta L}{\delta w_2} = \frac{1}{3} (0.5 + 1.5 - 1) = \frac{0.5}{3} = 0$$

$$\therefore \frac{\delta L}{\delta b} = \frac{1}{3} (-0.5 + 0.5 - 0.5) = \frac{-0.5}{3} = -0.167$$

$$\text{Updates} \rightarrow w_1 = 0 - 0.5 \times (-0.5) = 0.25$$

$$w_2 = 0 - 0.5 \times 0 = 0$$

$$b = 0 - 0.5 \times (-0.167) = 0.835$$