

Machine Learning

Exercise Sheet 2 Solution

Alper Hamza Ari
5771973

Chenqi Hao
5781375

Danil Skokov
5779466

Said Orfan Haidari
5781295

November 1, 2023

1 Linear Regression

Please see `question1.ipynb` program for the solution to the given problem.

2 Logistic Regression

2.1 Deriving update rule for gradient descent

Given the objective loss function (negative loglikelihood or binary cross entropy loss):

$$J = - \sum_{i=1}^N y_i \log(p_i) + (1 - y_i) \log(1 - p_i) \quad (1)$$

And predicted value of i th observation:

$$p_i = h_w(x_i) = P(y = 1 | X = x_i; w) = \sigma(x_i) \quad (2)$$

We want to derive update rule of gradient descent optimization for logistic regression, which aims to find the optimal weights to minimize objective loss.

Given update rule and its parameters:

$$w^{(t+1)} = w^{(t)} - \eta \frac{\partial J(D, w)}{\partial w} \quad (3)$$

where $t \in \mathbb{N}^+$

and given dataset $D = \{(x_i, y_i)\}_{i=1}^N$

For simplicity, i will remove i from all equations and i will perform all the operations in matrix form. Re-expressing what we had:

$$J = -y \log(p) - (1 - y) \log(1 - p) \quad (4)$$
$$p = \sigma(z) = \frac{1}{1 + e^{-z}}$$
$$z = w^\top x + b$$

By chain rule:

$$\begin{aligned} \frac{\partial J}{\partial w} &= \frac{\partial J}{\partial p} \cdot \frac{\partial p}{\partial w} \\ &= \frac{\partial J}{\partial p} \cdot \frac{\partial p}{\partial z} \cdot \frac{\partial z}{\partial w} \end{aligned} \quad (5)$$

And by calculating each partial derivative:

$$\begin{aligned}
\frac{\partial J}{\partial p} &= \frac{-y}{p} + \frac{1-y}{1-p} \\
\frac{\partial p}{\partial z} &= \frac{e^{-z}}{(1+e^{-1})^2} \\
\frac{\partial z}{\partial w} &= x
\end{aligned} \tag{6}$$

By putting pieces together:

$$\frac{\partial J}{\partial w} = \left(\frac{-y}{p} + \frac{1-y}{1-p} \right) \left(\frac{-e^{-z}}{(1+e^{-z})^2} \right) \cdot (x) \tag{7}$$

And recall that:

$$\begin{aligned}
p &= \frac{1}{1+e^{-z}} \\
e^{-z} &= \frac{1-p}{p}
\end{aligned} \tag{8}$$

Then our final partial derivative can be expressed again as follows:

$$\begin{aligned}
\frac{\partial J}{\partial w} &= \left(\frac{-y}{p} + \frac{1-y}{1-p} \right) \cdot \left(\frac{(1-p)p^2}{p} \right) \cdot (x) \\
&= (p-y)x
\end{aligned} \tag{9}$$

Which can also be expressed in different ways as follows:

$$\begin{aligned}
\frac{\partial J}{\partial w} &= \sum_i^N -(y_i - h_w(x_i)) \cdot x_i \\
&= -X^\top (y - p)
\end{aligned} \tag{10}$$

2.2 One step of gradient descent algorithm

We will append bias term to the end of weight array and define a new parameter array:

$$w = [w_1 \quad w_2 \quad b] \tag{11}$$

which was initially given as:

$$w = [0 \quad 0 \quad 0] \tag{12}$$

Then we will find predicted values for each observation x_n :

$$\begin{aligned}
&\text{Predicted value: } p = \sigma(w^\top x) \\
p_1(\text{or } \hat{y}_1) &= \sigma((w^\top x_1)) = \sigma(0) = 0.5 \\
p_2(\text{or } \hat{y}_2) &= \sigma((w^\top x_2)) = \sigma(0) = 0.5 \\
p_3(\text{or } \hat{y}_3) &= \sigma((w^\top x_3)) = \sigma(0) = 0.5 \\
p_4(\text{or } \hat{y}_4) &= \sigma((w^\top x_4)) = \sigma(0) = 0.5 \\
&\text{which gives us the } \hat{y} \text{ array} \\
\hat{y} &= [0.5 \quad 0.5 \quad 0.5 \quad 0.5]
\end{aligned} \tag{13}$$

Then we will find gradient vector:

$$\frac{\partial J}{\partial w} = \begin{bmatrix} \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_2} \\ \frac{\partial J}{\partial b} \end{bmatrix} = \begin{bmatrix} (\hat{y} - y)x_1 \\ (\hat{y} - y)x_2 \\ 0 \end{bmatrix}$$

where x_n represents array of observations for n th feature

$$x_1 = \begin{bmatrix} 2 \\ 3 \\ -4 \\ -2 \end{bmatrix}$$

$$x_2 = \begin{bmatrix} 4 \\ 3 \\ -2 \\ -6 \end{bmatrix}$$

y represents array of results

and \hat{y} represents array of predicted values

$$y = [1 \quad 1 \quad 0 \quad 0]$$

$$\hat{y} = [0.5 \quad 0.5 \quad 0.5 \quad 0.5]$$

Then calculating gradient vector

$$\frac{\partial J}{\partial w_1} = [-0.5 \quad -0.5 \quad 0.5 \quad 0.5] \cdot \begin{bmatrix} 2 \\ 3 \\ -4 \\ -2 \end{bmatrix} = -5.5$$

$$\frac{\partial J}{\partial w_2} = [-0.5 \quad -0.5 \quad 0.5 \quad 0.5] \cdot \begin{bmatrix} 4 \\ 3 \\ -2 \\ -6 \end{bmatrix} = -7.5$$

$$\frac{\partial J}{\partial w} = \begin{bmatrix} \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_2} \\ \frac{\partial J}{\partial b} \end{bmatrix} = \begin{bmatrix} -5.5 \\ -7.5 \\ 0 \end{bmatrix}$$

Also notice that $\frac{\partial J}{\partial b} = 0$ since we neglect bias in the formula $z = w^T x + b$ which leaves us with $z = w^T x$. Please check equation 4 above.

Calculating new weights after 1 step:

$$w^{(1)} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} - (0.5) \begin{bmatrix} -5.5 \\ -7.5 \\ 0 \end{bmatrix} = \begin{bmatrix} 2.75 \\ 3.75 \\ 0 \end{bmatrix} \quad (15)$$

Predicting $P(y = 1|X = [-1, 1])$:

$$\begin{aligned} \sigma(w^T x + b) &= \sigma(w_1 x_1 + w_2 x_2 + b) \\ &= \sigma(-2.75 + 3.75 + 0) \\ &= \sigma(1) \\ &= \frac{1}{1 + e^{-1}} \\ &= 0.731 \end{aligned} \quad (16)$$