

FOML Assignment1

SADINENI ABHINAY - CS21BTECH11055

NIMMALA AVINASH - CS21BTECH11039

Group 64

October 2023

Question 4

Let us consider the following notations

(i) $h_w(x)$ be model function

(ii) w be parameter

$$h_w(x) = g(w^T x) = \frac{1}{1 + e^{-w^T x}} \quad (1)$$

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

Let's assume that

$$P(y = 1|x; w) = h_w(x) \quad (3)$$

$$P(y = 0|x; w) = 1 - h_w(x) \quad (4)$$

$$P(y|x; w) = (h_w(x))^y (1 - h_w(x))^{(1-y)} \quad (5)$$

(a) (i) **Log likelihood:**

$$L(w) = p(y|x; w) \quad (6)$$

$$= \prod_{i=1}^n P(y_i|x_i; w) \quad (7)$$

$$= \prod_{i=1}^n (h_w(x_i))^{y_i} (1 - h_w(x_i))^{(1-y_i)} \quad (8)$$

$$l(w) = \log(L(w)) \quad (9)$$

$$= \sum_{i=1}^n [y_i \log(h_w(x_i)) + (1 - y_i) \log(1 - h_w(x_i))] \quad (10)$$

(ii) **Gradient:**

$$\frac{\partial l(w)}{\partial w_j} = \sum_{i=1}^n \left[y_i \frac{1}{h_w(x_i)} \frac{\partial}{\partial w_j} h_w(x_i) \right] \quad (11)$$

$$- \left[(1 - y_i) \frac{1}{1 - h_w(x_i)} \frac{\partial}{\partial w_j} h_w(x_i) \right] \quad (12)$$

$$\frac{\partial}{\partial w_j} h_w(x_i) = h_w(x_i) (1 - h_w(x_i)) \frac{\partial}{\partial w_j} (w^T x_i) \quad (13)$$

$$\frac{\partial l(w)}{\partial w_j} = \sum_{i=1}^n [y_i x_{ij} (1 - h_w(x_i)) - (1 - y_i) x_{ij} h_w(x_i)] \quad (14)$$

$$= \sum_{i=1}^n [x_{ij} (y_i - h_w(x_i))] \quad (15)$$

$$\nabla_w l(w) = \left[\frac{\partial l(w)}{\partial w_1} \frac{\partial l(w)}{\partial w_2} \dots \frac{\partial l(w)}{\partial w_n} \right]^T \quad (16)$$

$$= X^T (YV - HV) \quad (17)$$

Here X is design matrix of $[x \dots]$, YV is vector containing target values from data point 1 to n, HV is vector containing predicted values from data point 1 to n

(iii) **Hessian Matrix:** $(H_{l(w)})$

$$H_{ij} = \frac{\partial^2 l(w)}{\partial w_i \partial w_j} \quad (18)$$

$$= \frac{\partial}{\partial w_i} \left(\sum_{k=1}^n [x_{kj} (y_k - h_w(x_k))] \right) \quad (19)$$

$$= \sum_{k=1}^n [-x_{kj} x_{ki} h_w(x_k) (1 - h_w(x_k))] \quad (20)$$

$$= -X^T D X \quad (21)$$

where D is diagonal matrix with $D_{ii} = h_w(x_i)(1 - h_w(x_i))$

(iv) **Update equation:**

$$w := w + H_{l(w)}^{-1} \nabla_w l(w) \quad (22)$$

(v) **Algorithm:**

Algorithm 1: Gradient function

Data: w, x, y

Result: The gradient vector

1 **foreach** j **in** 1 to n **do**

2 $\text{partial_derivative}_j \leftarrow \sum_{i=1}^n x_{ij} (y_i - h_w(x_i))$

3 **end**

4 **return** $[\text{partial_derivative}_1, \dots, \text{partial_derivative}_n]^T$

Algorithm 2: Hessian matrix**Data:** w, x, y **Result:** The Hessian matrix of the function at w

```

1 Initialize an empty square matrix  $H$  with dimensions  $n \times n$ ;
2 foreach  $i$  in 1 to  $n$  do
3   foreach  $j$  in 1 to  $n$  do
4     // Compute each element of the Hessian matrix
4     //  $x_{ij}$  is the  $j$ -th component of  $x_i$ 
4      $H_{ij} \leftarrow \sum_{k=1}^n [-x_{kj} x_{ki} h_w(x_k) (1 - h_w(x_k))];$ 
5   end
6 end
7 return  $H$ ;

```

Algorithm 3: Newton-Raphson iterative method**Data:** initial_guess, x, y , $constant_1$, $constant_2$

```

1 iterations  $\leftarrow constant_1$ 
2 convergence  $\leftarrow constant_2$ 
3  $w \leftarrow initial\_guess$ 
4  $i \leftarrow 1$ 
5  $gradient_{new} \leftarrow gradient(w, x, y)$ 
6  $H \leftarrow hessian(w, x, y)$ 
7 while  $gradient_{new} > convergence$  and  $i \leq iterations$  do
8    $w^{i+1} \leftarrow w^i + inverse(H) \cdot gradient_{new}$ 
9    $i \leftarrow i + 1$ 
10   $gradient_{new} \leftarrow gradient(w^i, x, y)$ 
11   $H \leftarrow hessian(w^i, x, y)$ 
12 end
13 return  $w^i$ 

```

(b) Let us start by evaluating the updated parameter for each iteration.

$$w^{i+1} = w^i - inverse(H) \cdot gradient_{new} \quad (23)$$

$$= w^i - (X^T DX)^{-1} X^T (YV - HV) \quad (24)$$

$$= (X^T DX)^{-1} (X^T DX) w^i - (X^T DX)^{-1} X^T (YV - HV) \quad (25)$$

$$= (X^T DX)^{-1} ((X^T DX) w^i - X^T (YV - HV)) \quad (26)$$

$$= (X^T DX)^{-1} (X^T DP) \quad (27)$$

where $P = Xw^i - D^{-1}(YV - HV)$

D is invertible matrix since diagonals are non-zero since $w^T x_i$

- (i) The solution is similar to solution for **weighted least squares** $((\Phi^T R \Phi)^{-1} (\Phi^T R t))$ **in problem 3c**
- (ii) This solution for logistic regression has weights for its error function which are dependent on parameter and data.

- (iii) these weights are changed for every iteration D matrix changes as w changes
 - (iv) these are the reason it is called re weighted iterative method.
- (c) Error function of a logistic regression is given by

$$Err(h_w, y) = -l(w) \quad (28)$$

Error function is convex if its hessian matrix is semi-definite positive

$$H_{Err(w)} = -H_{l(w)} \quad (29)$$

$$= X^T D X \quad (30)$$

$$w H_{Err(w)} w^T = w X^T D X w^T \quad (31)$$

$$= (w X^T D^{1/2})(w X^T D^{1/2})^T \quad (32)$$

$$= ||w X^T D^{1/2}||^2 \quad (33)$$

Norm of vector / matrix is always greater than or equal to zero So Hessian matrix is positive semi-definite .therefore Error function is convex function on w.