

Name: Mayuresh R. Dindorkar  
Roll NO.: CS23 MTECH 14007

Name: Sanyam Kaul.  
Roll NO.: CS23 MTECH 14011

Q.4) For logistic regression, there is no longer a closed form solution due to nonlinearity of logistic sigmoid function. The error function can be minimized by an efficient iterative technique based on Newton-Rapson iterative optimization scheme.

Q) Provide the expressions of gradient, Hessian and update equations for Newton-Rapson optimization technique used to obtain the parameters in logistic regression model. Provide the algorithm, describing the methodology.

In order to obtain the expressions of Hessian, gradient & update equations, we need the error function  $E(\omega)$  of logistic regression. Hence, first we will find the error function for logistic regression.

### \* Logistic Regression:-

Consider the dataset  $\{\phi_n, t_n\}$ , where  $t_n \in \{0, 1\}$  and  $\phi_n = \phi(x_n)$  with  $n = 1, 2, \dots, N$ .

$\omega$  = parameter vector.

In logistic regression,

we model the probability of positive class ( $C_1$ ) as, of positive class ( $C_1$ )

$$P(C_1 | \phi) = y(\phi) = \sigma(\omega^T \phi) = \frac{1}{1 + e^{-\omega^T \phi}}$$

$$\therefore P(C_2 | \phi) = 1 - P(C_1 | \phi) = \frac{1}{1 + e^{\omega^T \phi}}$$

The likelihood function for logistic regression, can be written as,

$$P(\vec{y}|\vec{w}) = \prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n}$$

where,  $\vec{t} = (t_1, t_2, \dots, t_N)^T$  and  $y_n = p(c_1|\vec{w}_n)$ .

To find the error function  $E(w)$ , we take negative logarithm of the likelihood as below:-

$$E(w) = -\log P(\vec{t}|w)$$

$$E(w) = -\sum_{n=1}^N t_n \log y_n + (1-t_n) \log (1-y_n) \quad \leftarrow ①$$

which is also called "cross entropy error" function,

where  $y_n = g(w_n)$  and  $w_n = w^T \vec{x}_n$ .

### (I) Equation for gradient (Denoted by $\nabla E(w)$ )

The equation for gradient is obtained by taking derivative of  $E(w)$  w.r.t.  $w$ ,

$$\therefore \nabla E(w) = \frac{\partial}{\partial w} \sum_{n=1}^N (t_n - y_n) \vec{x}_n$$

$$\nabla E(w) = \vec{\Phi}^T (\vec{y} - \vec{t}) \quad \leftarrow ②$$

where  $\vec{\Phi}$  is "design matrix" of size  $N \times M$ ,

$$\vec{\Phi} \in \mathbb{R}^{N \times M} \quad \text{(any)} \quad = (\Phi_{ij} = (x_{ij}))_{N \times M}$$

$$\vec{y} = [y_1, y_2, \dots, y_N]$$

Equation ② represents the equation of gradient.

(II) Equation for Hessian:- (Denoted by  $H$ ). (III)

Hessian matrix  $H$  is obtained by taking double derivative of  $E(\omega)$ , viz if  $\omega$  is minimum then minimization

$$\therefore H = \nabla \nabla E(\omega) \quad \text{where } \omega = \omega$$

$$\frac{\partial}{\partial \omega} E(\omega) = 0 \quad \text{for minimum of } E(\omega)$$

$$\begin{aligned} &= \frac{\partial}{\partial \omega} \left[ \sum_{n=1}^N (y_n - t_n) \phi(x_n) \right] \quad \text{from ②} \\ &= \sum_{n=1}^N \phi(x_n) \left( \frac{\partial y_n}{\partial \omega} \right) \end{aligned}$$

$$\therefore H = \sum_{n=1}^N y_n (1-y_n) \Phi_n \Phi_n^T$$

$$\therefore H = \Phi^T R \Phi \quad \text{..... ③}$$

We have introduced the diagonal matrix  $R$  of dimensions  $N \times N$  with elements,

$$R_{nn} = y_n (1-y_n)$$

Equation ③ represents the equation of Hessian.

$$((I - \Phi)^T R - \omega \Phi^T \Phi)^T ((I - \Phi) =$$

$$((I - \Phi)^T R - \omega \Phi^T \Phi)^T = \omega$$

$$(I - \Phi)^T R - \omega \Phi^T \Phi = \omega$$

(III)

### Equation for update:

The update equation according to Newton Rapsori for minimizing error function  $E(\omega)$  is given by

$$\omega^{\text{new}} = \omega^{\text{old}} + \Delta \omega$$

where  $\Delta \omega$  is called 'stepsize' i.e.  $\Delta \omega = \omega^{\text{new}} - \omega^{\text{old}}$

We want to minimize the objective function  $\nabla E(\omega)$ .

We can expand  $E(\omega^{\text{new}})$  using Taylor Series around current estimate  $\omega^{\text{old}}$  as:

$$E(\omega^{\text{new}}) \approx E(\omega^{\text{old}}) + (\omega - \omega^{\text{old}})^T \nabla E(\omega^{\text{old}}) + \frac{1}{2} (\omega - \omega^{\text{old}})^T \cdot H \cdot (\omega - \omega^{\text{old}}).$$

where,

$\nabla E(\omega)$  is gradient vector and  $H$  is Hessian matrix.

Now, we differentiate it and equate to zero, to obtain  $\omega^{\text{new}}$  as:

$$\frac{\partial}{\partial \omega} \left[ E(\omega^{\text{old}}) + (\omega - \omega^{\text{old}})^T \cdot \nabla E(\omega^{\text{old}}) + \frac{1}{2} (\omega - \omega^{\text{old}})^T \cdot H \cdot (\omega - \omega^{\text{old}}) \right] = 0$$

Solving this we get,  $\boxed{\omega^{\text{new}} = \omega^{\text{old}} - H^{-1} \cdot \nabla E(\omega^{\text{old}})} \quad \dots \text{eqn ①}$

which represents the 'update equation' for newton rapsori technique.

After substituting values of  $H$  &  $\nabla E$ , the update equation becomes,

$$\omega^{\text{new}} = \omega^{\text{old}} - (\Phi^T R \Phi)^{-1} \Phi^T (\vec{y} - \vec{E})$$

$$= (\Phi^T R \Phi)^{-1} [ \Phi^T R \Phi \cdot \omega^{\text{old}} - \Phi^T (\vec{y} - \vec{E}) ]$$

$$\omega^{\text{new}} = (\Phi^T R \Phi)^{-1} \Phi^T R z$$

$$\text{where } z = \Phi^T \omega^{\text{old}} - R^{-1} (\vec{y} - \vec{E}).$$

$z$  is vector with  $N$  dimensions.

Chap 9  
Ex 9

### (\*) Algorithm:-

- ① Initialize  $\omega^0$ . with some random weights. (i)
- ② For each  $\tau = 1, 2, 3, \dots, N$  ( $N = \text{max. No. of iterations}$ )
  - (i) Compute  $\nabla E(\omega^\tau)$ , A gradient. (ii)
  - (ii) Compute H (Hessian matrix) using  $\omega^\tau$ . (iii)
  - (iii) Determine  $\omega^{\tau+1}$  using Newton-Rapson's update expression.  

$$\omega^{\tau+1} = \omega^\tau - H^{-1} \cdot \nabla E(\omega^\tau)$$
 (iv)
  - (iv) Check whether convergence criteria met or not.  
 If met, stop & come out of loop.
- ③ Return the estimated parameter vector  $\omega^*$ .

(Q.4) (b) Show that the Newton-Rapson update scheme is related to the weighted least squares problem described in the question 3(c). and explain why it is called the iterative reweighted least squares method.

In logistic regression, the Newton-Rapson update technique updates the parameter vector  $w$  using inverse of  $H$  i.e Hessian matrix. This update operation resembles a weighted least squares problem.

The update equation for logistic regression using Newton-Rapson is given by,

$$w^{\text{new}} = w^{\text{old}} - H^{-1} \nabla E(w^{\text{old}}) \quad \text{--- (1)}$$

$$\therefore w^{\text{new}} = (\Phi^T R \Phi)^{-1} \Phi^T R \cdot z \quad \text{--- (2)}$$

where,

$z$  is a vector with  $N$  dimensions.

$R$  = Diagonal matrix with  $N \times N$  dimensions with elements

$$R_{nn} = y_n(1-y_n).$$

$\Phi$  = Design matrix of size  $N \times M$ .

In (1), the  $H^{-1}$  acts as a weight matrix where the weight are determined by curvature of error function.

specifically, In particular, data points with steeper sigmoid function gradients contributes less to the update.

The weight matrix  $H^{-1}$  resembles  $R$  having elements  $r_{nn}$  as

$$r_{nn} = \delta(y_n) [1 - \delta(y_n)],$$

where,

$$y_n = \Phi^T x_n; \text{ in context of WLS method.}$$

Equation (2) closely resembles the form of normal equation of weighted least squares (WLS) problem.

As the matrix  $R$  is dependant on parameter vector  $w$ , the normal equation needs to be applied iteratively, using new weight vector  $w$  at each iteration to obtain the revised weight matrix  $R$ .

Therefore, this algorithm is often referred as Iterative Re-weighted least squares' (IRLS) method.

$$(1) \quad ({}^{(1)}w)^T H - b = {}^{(1)}w$$

$$(2) \quad S^{-1} E^T (E S^{-1} E) = w^T (w)$$

Step 1: Set initial value of  $w$  as  ${}^{(1)}w$ .  
Step 2: Compute  $S = E^T E$

Step 3: Compute  $b = E^T b$   
Step 4: Compute  $w = S^{-1} b$

$$\text{Step 5: Compute } [G(\beta) \delta - 1] G(\beta) \delta = \alpha$$

Step 6: If  $\alpha < \epsilon$  then stop, else go to step 4.

(Q.4) (c) Show that the error function of logistic regression is a convex function of  $w$  and has unique minimum with help of Hessian matrix.

→ First we will prove positive definiteness of Hessian matrix  $H$  and then will prove that the solution is unique minimum. to not get multiple minima

### I) Positive definiteness of Hessian matrix ( $H$ ):

we will show that  $R(\{y_n\})$  is positive definite matrix.

Hence, as per definition of positive definiteness,

$$u^T R u = \sum_n u_n \cdot y_n (1-y_n) \cdot u_n.$$

$$\therefore u^T R u = \sum_n u_n^2 \cdot y_n (1-y_n)$$

where  $u$  is any arbitrary vector and  $R$  is design matrix.

for  $y_n$  = output of logistic regression model for input  $x_n$ .

This equation will be positive for any arbitrary non zero vector  $u$ , because sigmoid ensures  $0 < y_n < 1$ .

By using same analogy,

Consider a Hessian matrix  $H$  & arbitrary vector  $u$ ,

$$\therefore u^T H u = u^T \Phi^T R \Phi u$$

$$\text{where } \Phi = \text{design matrix.}$$

$$\therefore u^T H u = v^T H v \quad \dots v = u^T \Phi^T.$$

$$= \sum_n v_n^2 \cdot y_n (1-y_n) > 0.$$

Hence this shows that hessian matrix  $H$  is also positive definite. Hence, error function of logistic regression is a convex function.

11

Unique minimum: visiting every node once with initial cost = 0

After expanding the Taylor series around minimum point  $w^*$ , we get our objective function as:

$$E(\omega^* + \Delta\omega) = E(\omega^*) + \Delta\omega E(\omega^*) \Delta\omega + (\Delta\omega)^T \cdot H(\omega) \cdot (\Delta\omega) + Z(||\Delta\omega||^3)$$

where,  $H$  is function of  $w$  because all  $y_n$ 's are function of  $w$ .

At minimum, the first order term  $\Delta E(\omega^*) \cdot \Delta\omega = 0$  because derivative is zero at minimum point.

Similarly, the higher order terms denoted by  $2(\|\Delta w\|)^3$  are negligible.

The remaining quadratic term  $(\Delta w)^T H (\Delta w)$  gets minimized when  $\Delta w = 0$ , hence our objective function also gets minimized at  $\Delta w = 0$ .

Hence, by using properties of positive definiteness of Taylor series expansion we have shown that  $E(w)$  has unique minimum at  $w^*$ .