

## Logistic Regression

linear regression —

$$y = \underline{x}^T \underline{\theta} + \epsilon$$

logistic function is —

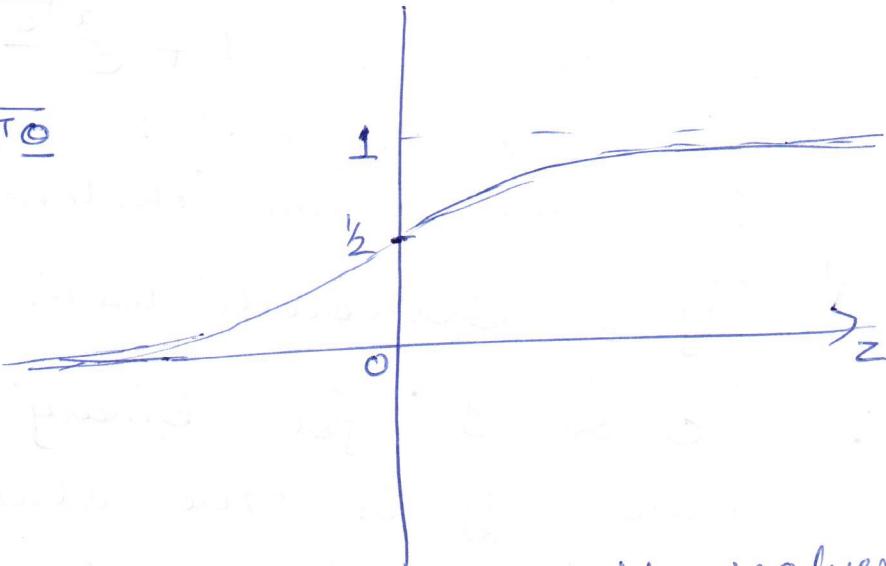
$$y = \frac{1}{1 + e^{-\underline{x}^T \underline{\theta}}}$$

$$\text{let } \underline{x}^T \underline{\theta} = z$$

$$z=0, \quad y = \frac{1}{2}$$

$$z=\infty, \quad y = 1$$

$$z=-\infty, \quad y = 0$$



$\therefore y$  looks like a prob. fn. with values between 0 & 1

Also, this function  $y$  is smooth (and hence differentiable), & also symmetric.

For binary classification,

$y \geq \frac{1}{2} \rightarrow$  predict class-1

$y < \frac{1}{2} \rightarrow$  " class-0.

$$\therefore P(y=1/\underline{x}) = \text{logistic}(\underline{x}, \underline{\theta})$$

$$\begin{aligned}
 &= \frac{1}{1 + e^{-\underline{x}^T \underline{\theta}}} \\
 &= \frac{e^{\underline{x}^T \underline{\theta}}}{1 + e^{\underline{x}^T \underline{\theta}}}
 \end{aligned}$$

$$\begin{aligned}\therefore p(y=0/x) &= 1 - p(y=1/x) \\ &= 1 - \frac{1}{1 + e^{-x^T \theta}} = \frac{e^{-x^T \theta}}{1 + e^{-x^T \theta}} \\ &= \frac{1}{1 + e^{x^T \theta}}\end{aligned}$$

Use maximum likelihood estimator —

Actual  
 $y$  = Bernoulli with two values  
 0 & 1 for binary classification problem,  
 where  $y$  is the class label.  
 Let there be  $n$  no. of independent  
 observations of  $y$  for  $i = 1, 2, \dots, n$   
 generated by  $n$  feature vectors  
 $x_1, \dots, x_n$ .

Conditional likelihood fn is —

$$p(y/x, \theta) = \prod_{i=1}^n p(y_i/x_i, \theta) = l(\theta)$$

max. the likelihood w.r.t  $\theta$  will provide its estimate

$$\hat{\theta} = \max_{\theta} \prod_{i=1}^n p(y_i/x_i, \theta)$$

$$= \max_{\theta} \prod_{i=1}^n \left( p(y_i=1/x_i, \theta) \right)^{y_i} \left( p(y_i=0/x_i, \theta) \right)^{1-y_i}$$

Define log likelihood =  $L(\underline{\theta}) = \ln(L(\underline{\theta}))$

maximizing the  $\ln(L(\underline{\theta}))$  will give the same result as of maximizing  $L(\underline{\theta})$  w.r.t.  $\underline{\theta}$  because  $\ln$  is a monotonic fn. -

$$\therefore \hat{\underline{\theta}} = \max_{\underline{\theta}} L(\underline{\theta})$$

$$= \max_{\underline{\theta}} \left\{ \sum_{i=1}^n y_i \ln(p(y_i=1/x_i, \underline{\theta})) \right.$$

$$\left. + \sum_{i=1}^n (1-y_i) \ln(p(y_i=0/x_i, \underline{\theta})) \right\}$$

Define loss fn. as -

$$J(\underline{\theta}) = -L(\underline{\theta})$$

maximizing  $J(\underline{\theta})$   $\equiv$  minimizing  $J(\underline{\theta})$

$$\therefore \hat{\underline{\theta}} = \min_{\underline{\theta}} J(\underline{\theta}) = \min_{\underline{\theta}} -L(\underline{\theta})$$

$$= \min_{\underline{\theta}} \left\{ - \sum_{i=1}^n y_i \ln(p(y_i=1/x_i, \underline{\theta})) \right. \\ \left. - \sum_{i=1}^n (1-y_i) \ln(p(y_i=0/x_i, \underline{\theta})) \right\}$$

Binary cross Entropy Loss fn.  
(BCE)

$$\text{Entropy} = - \sum_i p_i \ln p_i$$

$$\text{cross entropy} = - \sum_i q_i \ln p_i$$

$$J(\underline{\theta}) = - \sum_{i=1}^n y_i \ln \left( \frac{e^{z_i}}{1+e^{z_i}} \right) - \sum_{i=1}^n (1-y_i) \ln \left( \frac{1}{1+e^{z_i}} \right)$$

where  $z_i = \underline{x}_i^T \underline{\theta}$

minimizes  $J(\underline{\theta})$  by taking the derivative of  $J(\underline{\theta})$  or computing the gradient of  $J(\underline{\theta})$ :

$$\begin{aligned} \frac{\partial J(\underline{\theta})}{\partial \underline{\theta}} &= - \frac{\partial}{\partial \underline{\theta}} \left[ \sum_{i=1}^n y_i \ln(e^{z_i}) + \sum_{i=1}^n (1-y_i) \ln(1+e^{z_i}) \right] \\ &= - \sum_{i=1}^n \left[ \frac{\partial}{\partial \underline{\theta}} \left( y_i \ln(e^{\underline{x}_i^T \underline{\theta}}) + (1-y_i) \ln(1+e^{\underline{x}_i^T \underline{\theta}}) \right) \right] \\ &= - \frac{\partial}{\partial \underline{\theta}} \left[ \sum_{i=1}^n \left( y_i \ln e^{\underline{x}_i^T \underline{\theta}} - \ln(1+e^{\underline{x}_i^T \underline{\theta}}) \right) \right] \\ &= - \sum_{i=1}^n \left( \frac{y_i \underline{x}_i e^{\underline{x}_i^T \underline{\theta}}}{e^{\underline{x}_i^T \underline{\theta}}} - \frac{\underline{x}_i e^{\underline{x}_i^T \underline{\theta}}}{1+e^{\underline{x}_i^T \underline{\theta}}} \right) \\ &= - \sum_{i=1}^n \underline{x}_i (y_i - p(y_i=1/\underline{x}_i^T \underline{\theta})) \end{aligned}$$

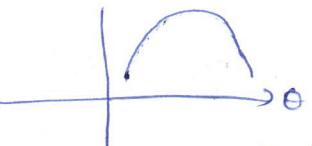
gradient looks like  $-\sum_{i=1}^n \underline{x}_i e_i$

if  $y_i = 1 \rightarrow$  we would like  $p(y_i=1/\underline{x}_i^T \underline{\theta}) \rightarrow 1$

if  $y_i = 0 \rightarrow$  " " " "  $p(y_i=1/\underline{x}_i^T \underline{\theta}) \rightarrow 0$

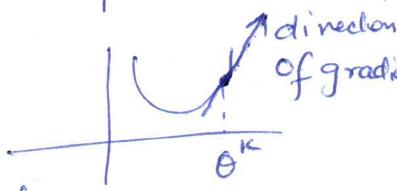
to keep the error small.

$L(\underline{\theta})$  is a concave fn. in  $\underline{\theta}$



3

&  $J(\underline{\theta})$  is " convex fn. in  $\underline{\theta}$



$\therefore \text{max. } L(\underline{\theta}) \equiv \text{gradient ascent}$

&  $\min. J(\underline{\theta}) \equiv$  gradient descent.

$J(\underline{\theta}) \equiv$  global min.  
will converge to

if we move iteratively as -

$$\underline{\theta}^{(k+1)} = \underline{\theta}^{(k)} - \eta \left( \frac{\partial J(\underline{\theta})}{\partial \underline{\theta}} \right)$$

- current value of  $\underline{\theta}$  in the

$\underline{\theta}^{(k)}$  =  $\underline{\theta}$  in the  
 $k^{\text{th}}$  iteration

$\underline{\theta}^{(k+1)}$  = value of  $\underline{\theta}$  in the  
 $(k+1)^{\text{th}}$  iteration

$\eta$  = learning rate

$\frac{\partial J(\underline{\theta})}{\partial \underline{\theta}}$  = move in the direction

opposite to the direction of gradient  
to reach to min.

$$\therefore \underline{\theta}^{(k+1)} = \underline{\theta}^{(k)} - \eta \frac{\partial J(\underline{\theta})}{\partial \underline{\theta}}$$

$$= \underline{\theta}^{(k)} + \eta \sum_{i=1}^n x_i e_i$$

Gradient  
Descent  
Method

$$\frac{\partial J(\underline{\theta})}{\partial \underline{\theta}} = - \sum_{i=1}^n \underline{x}_i \left( y_i - p(y_i=1/\underline{x}_i, \underline{\theta}) \right)$$

$$= - \underbrace{\underline{X}^T}_{p \times n} \left( \underline{y} - \underline{p} \right)_{n \times 1} = - \underline{X}^T \underline{e}$$

$$= - \left( \underline{x}_1 \underline{e}_1 + \underline{x}_2 \underline{e}_2 + \dots + \underline{x}_n \underline{e}_n \right)$$

Since

$$\underline{X} = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_n^T \end{bmatrix}$$

$\underline{\theta} \equiv p \times 1$   
vector

$$\therefore \underline{X}^T \underline{e} = \begin{bmatrix} \underline{x}_1 & \underline{x}_2 & \dots & \underline{x}_n \end{bmatrix}_{p \times n} \begin{bmatrix} \underline{e}_1 \\ \underline{e}_2 \\ \vdots \\ \underline{e}_n \end{bmatrix}$$

$$= \underline{x}_1 \underline{e}_1 + \dots + \underline{x}_n \underline{e}_n$$

$$\underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} ; \quad \underline{p} = \begin{bmatrix} p(y_1=1/\underline{x}_1, \underline{\theta}) \\ \vdots \\ p(y_n=1/\underline{x}_n, \underline{\theta}) \end{bmatrix}_{n \times 1}$$

$\frac{\partial^2 J(\underline{\theta})}{\partial \underline{\theta} \partial \underline{\theta}^T}$  Hessian matrix

$$= \frac{\partial}{\partial \underline{\theta}^T} \left( - \underline{X}^T (\underline{y} - \underline{p}) \right)$$

$p \times 1$  matrix

$p \times p$  matrix

$$\begin{aligned}
 \frac{\partial^2 J(\underline{\theta})}{\partial \underline{\theta} \partial \underline{\theta}^T} &= + \frac{\partial}{\partial \underline{\theta}^T} \left( -X^T (y - b) \right) \\
 &= \frac{\partial}{\partial \underline{\theta}^T} \left( -X^T y + X^T \left( \begin{array}{c} \frac{1}{1+e^{-x_1^T \underline{\theta}}} \\ \vdots \\ \frac{1}{1+e^{-x_n^T \underline{\theta}}} \end{array} \right) \right) \\
 &\quad \text{not a fn. of } \underline{\theta} \\
 &= \frac{\partial}{\partial \underline{\theta}^T} \left[ X^T \left( \begin{array}{c} \frac{1}{1+e^{-x_1^T \underline{\theta}}} \\ \vdots \\ \frac{1}{1+e^{-x_n^T \underline{\theta}}} \end{array} \right) \right] = \frac{\partial}{\partial \underline{\theta}^T} (X^T b) \\
 &= X^T \left[ \frac{\partial}{\partial \underline{\theta}^T} b \right]
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial \underline{\theta}^T} b &= \left[ \frac{\partial b}{\partial \theta_1}, \frac{\partial b}{\partial \theta_2}, \dots, \frac{\partial b}{\partial \theta_p} \right]^T \\
 \frac{\partial b}{\partial \theta_1} &= \frac{\partial}{\partial \theta_1} \left[ \begin{array}{c} \frac{1}{1+e^{-x_1^T \underline{\theta}}} \\ \vdots \\ \frac{1}{1+e^{-x_n^T \underline{\theta}}} \end{array} \right] = \left[ \begin{array}{c} x_{11} \frac{e^{-x_1^T \underline{\theta}}}{(1+e^{-x_1^T \underline{\theta}})^2} \\ \vdots \\ x_{n1} \frac{e^{-x_n^T \underline{\theta}}}{(1+e^{-x_n^T \underline{\theta}})^2} \end{array} \right] \\
 &= \left[ \begin{array}{c} x_{11} b(y_1=1/x_{11}, \underline{\theta}) (1 - b(y_1=1/x_{11}, \underline{\theta})) \\ \vdots \\ x_{n1} b(y_n=1/x_{n1}, \underline{\theta}) (1 - b(y_n=1/x_{n1}, \underline{\theta})) \end{array} \right]
 \end{aligned}$$

$$\frac{\partial \underline{p}}{\partial \theta_1} = \begin{bmatrix} x_{11} & p_1(1-p_1) \\ \vdots & \vdots \\ \vdots & \vdots \\ x_{n1} & p_n(1-p_n) \end{bmatrix}$$

$$= \begin{bmatrix} p_1(1-p_1) & 0 & \dots & 0 \\ 0 & p_2(1-p_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} x_{11} \\ \vdots \\ \vdots \\ x_{n1} \end{bmatrix}$$

$m \times n$  diagonal matrix

$$= W \begin{bmatrix} x_{11} \\ \vdots \\ \vdots \\ x_{n1} \end{bmatrix}$$

Similarly,

$$\frac{\partial \underline{p}}{\partial \theta_2} = W \begin{bmatrix} x_{12} \\ \vdots \\ \vdots \\ x_{n2} \end{bmatrix}$$

$$\frac{\partial \underline{p}}{\partial \theta_p} = W \begin{bmatrix} x_{1n} \\ \vdots \\ \vdots \\ x_{nn} \end{bmatrix}$$

$$\therefore \left[ \frac{\partial \underline{p}}{\partial \theta_1}, \frac{\partial \underline{p}}{\partial \theta_p} \right] = W \begin{bmatrix} x_{11} & \dots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nn} \end{bmatrix}$$

$$= W X$$

$$\frac{\partial^2 J(\underline{\theta})}{\partial \underline{\theta} \partial \underline{\theta}^T} = X^T \begin{bmatrix} \frac{\partial}{\partial \underline{\theta}} & b \end{bmatrix}$$

$$= X^T \text{Id} X$$

Newton's algorithm of optimization -

This can be used on a twice differentiable fn. & is based on the quadratic Taylor's expansion.

Let us consider the minimization of a fn. of one variable — say global min. is at  $\underline{\theta}^*$ .

Guess an 'initial'  $\underline{\theta} = \underline{\theta}^{(0)}$  ← can be zero.  
Use <sup>2nd order</sup> Taylor's Series expansion in the neighborhood of  $\underline{\theta}^{(0)}$  —

$$f(\underline{\theta}) \approx f(\underline{\theta}^{(0)}) + (\underline{\theta} - \underline{\theta}^{(0)}) \frac{df}{d\underline{\theta}} \Big|_{\underline{\theta}=\underline{\theta}^{(0)}}$$

$$+ \frac{1}{2} (\underline{\theta} - \underline{\theta}^{(0)})^2 \frac{d^2 f}{d \underline{\theta}^2} \Big|_{\underline{\theta}=\underline{\theta}^{(0)}}$$

$$= f(\underline{\theta}^{(0)}) + (\underline{\theta} - \underline{\theta}^{(0)}) f'(\underline{\theta}^{(0)}) + \frac{1}{2} (\underline{\theta} - \underline{\theta}^{(0)})^2 f''(\underline{\theta}^{(0)})$$

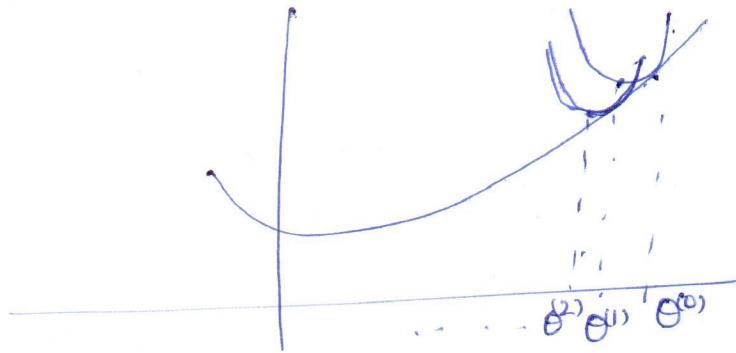
Take the derivative & equate to zero —

$$f'(\underline{\theta}) = 0 = f'(\underline{\theta}^{(0)}) + \frac{1}{2} 2(\underline{\theta} - \underline{\theta}^{(0)}) f''(\underline{\theta}^{(0)})$$

$$\Rightarrow \underline{\theta}^{(1)} - \underline{\theta}^{(0)} = - \frac{f'(\underline{\theta}^{(0)})}{f''(\underline{\theta}^{(0)})}$$

$$\Rightarrow \underline{\theta}^{(1)} = \underline{\theta}^{(0)} - \frac{f'(\underline{\theta}^{(0)})}{f''(\underline{\theta}^{(0)})}$$

$$\text{or } \underline{\theta}^{(1)} = \underline{\theta}^{(0)} - (\nabla^2 f(\underline{\theta}^{(0)}))^{-1} \nabla f(\underline{\theta}^{(0)})$$



Newton's method  
of moving to the bottom of  
the curve.

$\downarrow$        $x$

converges faster than the gradient  
descent method -  
if we apply this to our problem -

$$\begin{aligned}
 \underline{\theta}^{\text{new}} &= \underline{\theta}^{\text{old}} \\
 \underline{\theta}^{(1)} &= \underline{\theta}^{(0)} - (\nabla^2 J(\underline{\theta}))^{-1} \nabla (J(\underline{\theta})) \\
 &= \underline{\theta}^{(0)} - (x^T W^{(0)} x)^{-1} (-x^T (y - \underline{p}^{(0)})) \\
 &= \underline{\theta}^{(0)} + (x^T W^{(0)} x)^{-1} (x^T (y - \underline{p}^{(0)})) \\
 (x^T W^{(0)} x) \underline{\theta}^{(1)} &= (x^T W^{(0)} x) \underline{\theta}^{(0)} + x^T (y - \underline{p}^{(0)}) \\
 &= x^T W^{(0)} \left[ x \underline{\theta}^{(0)} + W^{(0)\top} (y - \underline{p}^{(0)}) \right] \\
 &= x^T W^{(0)} \underline{z}^{(0)}
 \end{aligned}$$

$$\begin{aligned}
 (\underline{x}^T \underline{w}^{(0)} \underline{x}) \underline{\theta}^{(0)} &= \underline{x}^T \underline{w}^{(0)} \underline{z}^{(0)} \\
 \underline{\theta}^{(1)} &= (\underline{x}^T \underline{w}^{(0)} \underline{x})^{-1} \underline{x}^T \underline{w}^{(0)} \underline{z}^{(0)} \\
 \underline{\theta}^{(k+1)} &= (\underline{x}^T \underline{w}^{(k)} \underline{x})^{-1} \underline{x}^T \underline{w}^{(k)} \underline{z}^{(k)} \quad \text{--- (1)}
 \end{aligned}$$

This update solution is also called as the Iteratively reweighted least squares (IRLS) method.

remember least squares soln -

$$\underline{\theta} = (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{y} \quad \text{--- (2)}$$

Compare (1) & (2) -

there is a diagonal weight matrix in (1)

$$\underline{w}^{(0)} = \left( \frac{e^{-\underline{x}_1^T \underline{\theta}^{(0)}}}{1 + e^{-\underline{x}_1^T \underline{\theta}^{(0)}}} \right)^2$$

$$\left( \frac{e^{-\underline{x}_n^T \underline{\theta}^{(0)}}}{1 + e^{-\underline{x}_n^T \underline{\theta}^{(0)}}} \right)^2$$

$$\underline{w}^{(k)} = \left( \frac{e^{-\underline{x}_1^T \underline{\theta}^{(k)}}}{1 + e^{-\underline{x}_1^T \underline{\theta}^{(k)}}} \right)^2$$

$$\left( \frac{e^{-\underline{x}_n^T \underline{\theta}^{(k)}}}{1 + e^{-\underline{x}_n^T \underline{\theta}^{(k)}}} \right)^2$$