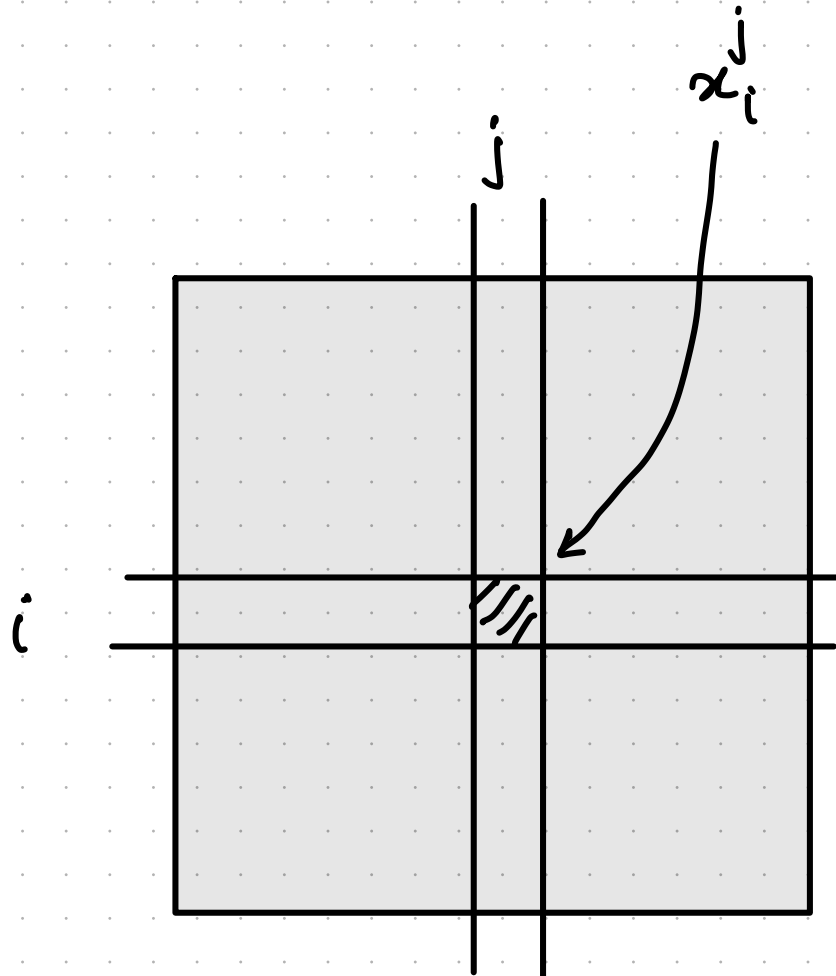


CONVEXITY OF CROSS ENTROPY LOSS

$$J(\theta) = - \sum_{i=1}^n (y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i))$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^n (\hat{y}_i - y_i) x_i^j$$

Where
 $x_i^j \leftarrow j^{\text{th}}$ element
 $x_i \leftarrow i^{\text{th}}$ data point



$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^n (\hat{y}_i - y_i) x_i^j$$

$$H = \begin{bmatrix} \frac{\partial}{\partial \theta_1} \frac{\partial J(\theta)}{\partial \theta_1} & \dots & \dots & \dots \\ \frac{\partial}{\partial \theta_2} \frac{\partial J(\theta)}{\partial \theta_1} & \dots & \dots & \dots \\ \vdots & \dots & \dots & \dots \\ \frac{\partial}{\partial \theta_2} \frac{\partial J(\theta)}{\partial \theta_2} & \dots & \dots & \dots \end{bmatrix}$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^N (\hat{y}_i - y_i) x_i^j$$

$$H = \begin{bmatrix} \frac{\partial}{\partial \theta_1} \frac{\partial J(\theta)}{\partial \theta_1} & & & \\ & \ddots & & \\ & & \frac{\partial}{\partial \theta_2} \frac{\partial J(\theta)}{\partial \theta_1} & \\ & & \vdots & \\ & & & \frac{\partial}{\partial \theta_d} \frac{\partial J(\theta)}{\partial \theta_d} \end{bmatrix}$$

let us compute a $\frac{\partial}{\partial \theta_j} \frac{\partial J(\theta)}{\partial \theta_k}$ as H_{jk} entry

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^N (\hat{y}_i - y_i) x_i^j$$

let us compute a $\frac{\partial}{\partial \theta_j} \frac{\partial J(\theta)}{\partial \theta_k}$ as H_{jk} entry

$$\frac{\partial J(\theta)}{\partial \theta_k} = \sum_{i=1}^N (\hat{y}_i - y_i) x_i^k$$

$$\frac{\partial}{\partial \theta_j} \frac{\partial J(\theta)}{\partial \theta_k} = \frac{\partial}{\partial \theta_j} \sum_{i=1}^N (\hat{y}_i - y_i) x_i^k$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^N (\hat{y}_i - y_i) x_i^j$$

let us compute a $\frac{\partial^2 J(\theta)}{\partial \theta_j \partial \theta_k}$ as H_{jk} entry

$$\frac{\partial J(\theta)}{\partial \theta_k} = \sum_{i=1}^N (\hat{y}_i - y_i) x_i^k$$

$$\frac{\partial}{\partial \theta_j} \frac{\partial J(\theta)}{\partial \theta_k} = \frac{\partial}{\partial \theta_j} \sum_{i=1}^N (\hat{y}_i - y_i) x_i^k$$

$$H_{jk} = \frac{\partial}{\partial \theta_j} \sum_{i=1}^N \hat{y}_i x_i^k$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^N (\hat{y}_i - y_i) x_i^j$$

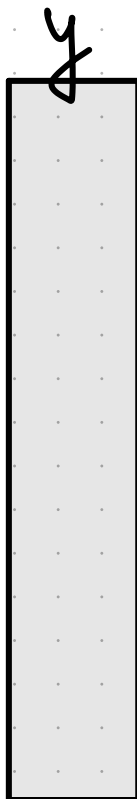
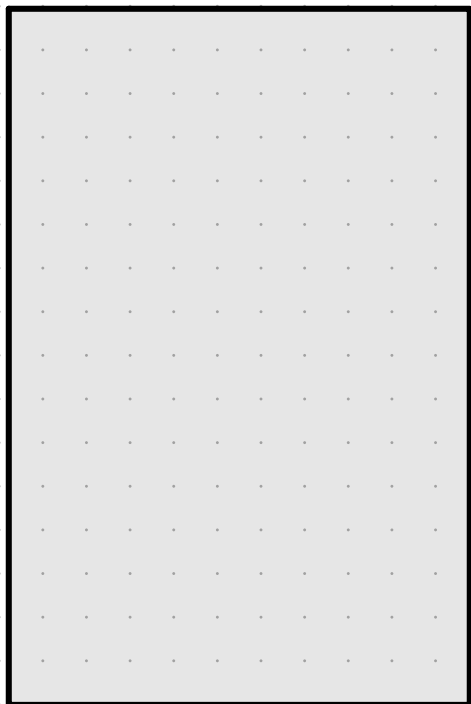
$$H_{jk} = \frac{\partial}{\partial \theta_j} \sum_{i=1}^N \hat{y}_i x_i^k$$

$$= \sum_{i=1}^N \hat{y}_i (1 - \hat{y}_i) x_i^k + \frac{\partial}{\partial \theta_j} (x_i^1 \theta_1 + x_i^2 \theta_2 + \dots)$$

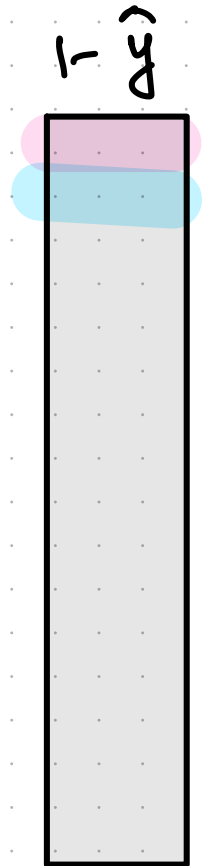
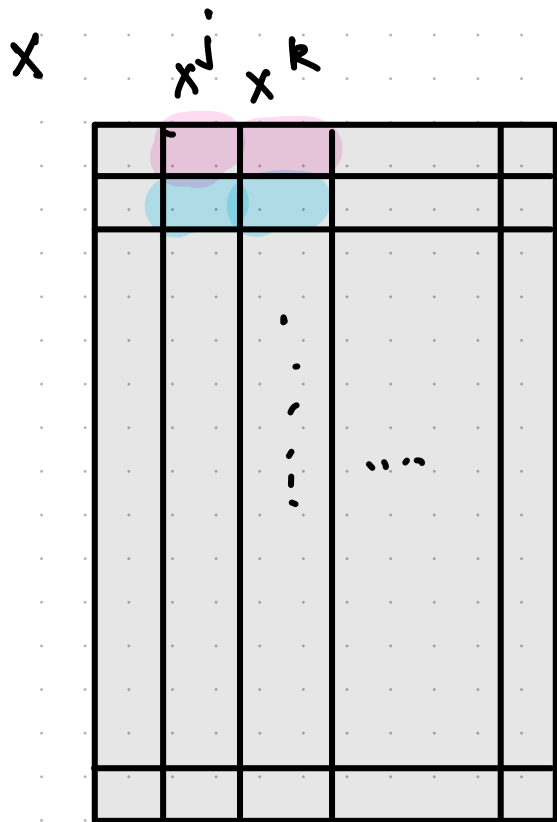
$$= \sum_{i=1}^N \hat{y}_i (1 - \hat{y}_i) x_i^k x_i^j$$

$$H_{jk} = \sum_{i=1}^N \hat{y}_i (1 - \hat{y}_i) x_i^k x_i^j$$

x

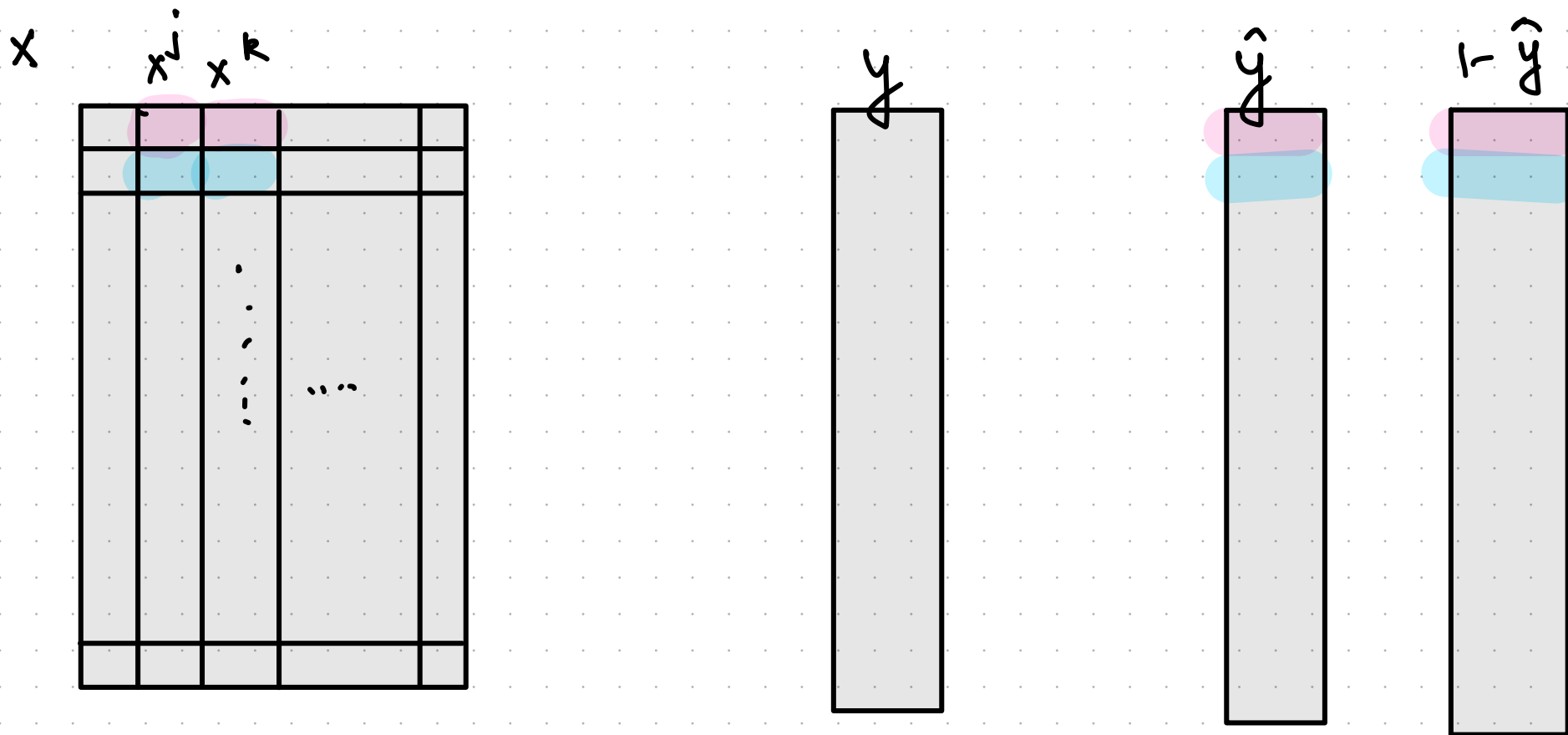


$$H_{jk} = \sum_{i=1}^N \hat{y}_i (1 - \hat{y}_i) x_i^k x_i^j$$



Say $j=2; k=3$

$$H_{jk} = \sum_{i=1}^N \hat{y}_i (1 - \hat{y}_i) x_i^k x_i^j$$



$$H_{jk} = x_j^T \begin{bmatrix} \hat{y}_1(1-\hat{y}_1) & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ \hat{y}_N(1-\hat{y}_N) & \dots & \dots & \dots \end{bmatrix} x^k_{N \times 1} = x_j^T \text{diag}(\hat{y} \odot (1 - \hat{y})) x^k$$

↑
Element-wise mult

$$H_{jk} = \sum_{i=1}^N \hat{y}_i (1 - \hat{y}_i) x_i^k x_i^j$$

$$H_{jk} = x_j^T \begin{matrix} 1 \times N \\ \left[\begin{array}{cccc} \hat{y}_1(1-\hat{y}_1) & \dots & 0 & 0 \\ & \ddots & & \\ & & \hat{y}_N(1-\hat{y}_N) & \end{array} \right]_{N \times N} \end{matrix} x^k \quad N \times 1$$

$$H_{jk} = x_j^T \text{diag}(\hat{y} \circ (1 - \hat{y})) x^k$$

$$H_{jk} = \sum_{i=1}^N \hat{y}_i (1 - \hat{y}_i) x_i^k x_i^j$$

$$H_{jk} = x_j^T \begin{bmatrix} \hat{y}_1(1-\hat{y}_1) & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ y_n(1-\hat{y}_n) \end{bmatrix} x^k$$

$1 \times N$ $N \times N$ $N \times 1$

$$H_{jk} = x_j^T \text{diag}(\hat{y}_i(1-\hat{y}_i)) x^k = X_j^T D x^k$$

What can we say about D ?

$$H_{jk} = \sum_{i=1}^N \hat{y}_i (1 - \hat{y}_i) x_i^k x_i^j$$

$$H_{jk} = x_i^T \begin{bmatrix} \hat{y}_1(1-\hat{y}_1) & \dots & 0 & 0 \\ \vdots & & & \\ \hat{y}_N(1-\hat{y}_N) \end{bmatrix} x^k$$

$1 \times N$ $N \times N$ $N \times 1$

$$H_{jk} = x_j^T \text{diag}(\hat{y}_i(1-\hat{y}_i)) x^k = X_j^T D x^k$$

What can we say about D ?

Each $\hat{y}_i \in [0, 1]$; $1 - \hat{y}_i \in [0, 1]$

$$\min \hat{y}_i [1 - \hat{y}_i] = 0.5 \times 0.5 = 0.25$$

$$H_{jk} = x_j^T \text{diag}(\hat{y} \odot (1 - \hat{y})) x^k = X_j^T D x^k$$

What can we say about D ?

$$\text{Each } \hat{y} \in [0, 1]; \quad 1 - \hat{y} \in [0, 1]$$

$$\text{Max } \hat{y} [1 - \hat{y}] = 0.5 \times 0.5 = 0.25$$

Thus D is of type

$$D = \begin{bmatrix} 0 \leq D_{11} \leq 0.25 & 0 & 0 & \dots \\ 0 & \cdot & & \\ 0 & & \cdot & \\ \vdots & & & \ddots \\ & & & & 0 \leq D_{nn} \leq 0.25 \end{bmatrix}$$

$$H_{jk} = x_j^T \text{diag}(\hat{y} \odot (1 - \hat{y})) x^k = X_j^T D x^k$$

What can we say about D ?

Each $\hat{y} \in [0, 1]$; $1 - \hat{y} \in [0, 1]$

$$\max \hat{y} [1 - \hat{y}] = 0.5 * 0.5 = 0.25$$

Thus D is diagonal matrix where diagonal

entries are b/w 0 and 0.25

$$H_{j^k} = x_j^T \text{diag}(\hat{y} \odot (r\hat{y})) x^k = X_j^T D x^k$$

What does H look like?

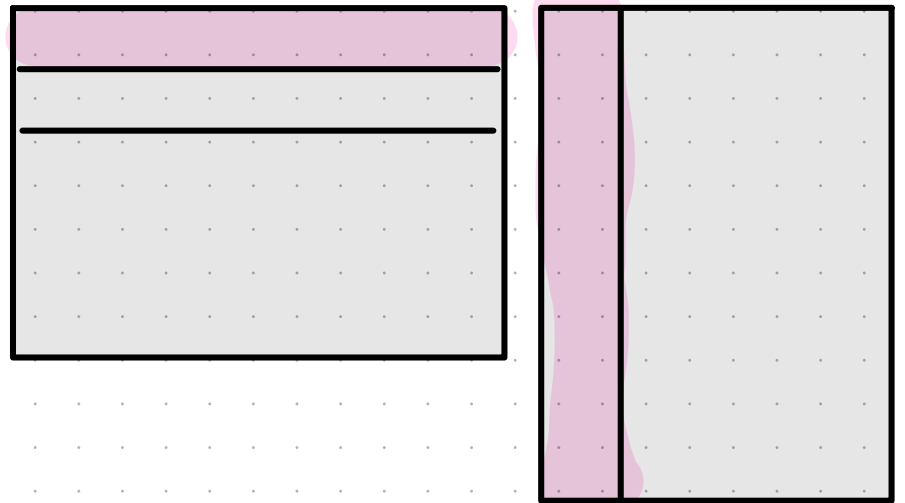
$$H_{jk} = x_j^T \text{diag}(\hat{y} \odot (r\hat{y})) x^k = x_j^T D x^k$$

What does H look like?

$$H_{11} = x_1^T D x_1$$

$$H_{21} = x_2^T D x_1$$

⋮



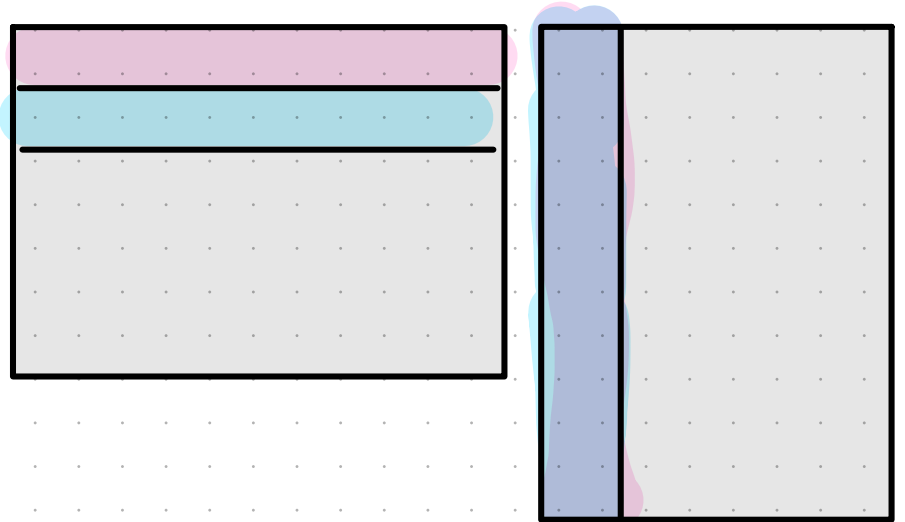
$$H_{jk} = x_j^T \text{diag}(\hat{y} \odot (r\hat{y})) x^k = x_j^T D x^k$$

What does H look like?

$$H_{11} = x^1{}^T D x^1$$

$$H_{21} = x^2{}^T D x^1$$

⋮



$$H_{jk} = x_j^T \text{diag}(\hat{y} \odot (r\hat{y})) x^k = x_j^T D x^k$$

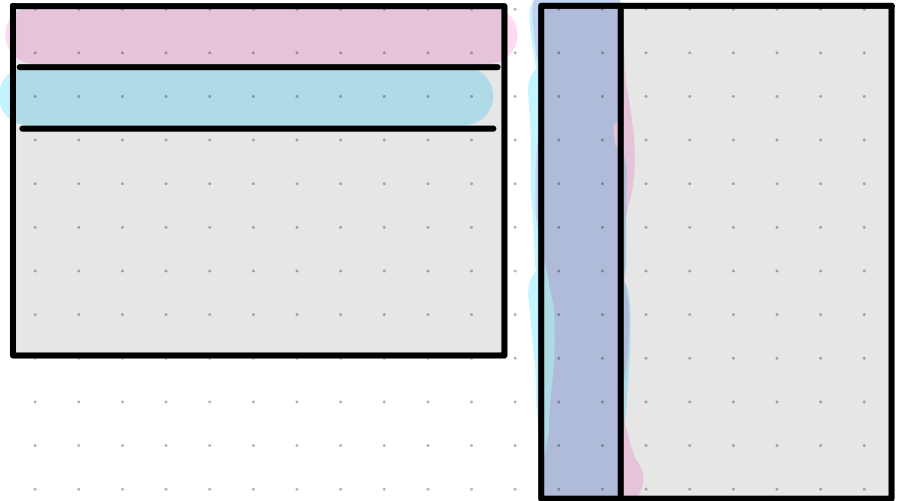
What does H look like?

$$H_{11} = x_1^T D x_1$$

$$H_{21} = x_2^T D x_1$$

⋮

$$H = X^T D X$$



$$H = X^T D X$$

Q: Prove H is P.S.D matrix

$$H = X^T D X$$

Q: Prove H is P.S.D matrix

$$v^T H v = v^T X^T D X v$$

$$\text{let } X v = z$$

$$v^T H v = z^T D z$$

$$= \sum_{i=1}^N d_{ii} z_i^2$$

where $0 \leq d_{ii} \leq 0.25$

$\geq 0 \quad \therefore H \text{ is P.S.D. } \Rightarrow J(\theta) \text{ is convex}$

Iteratively Reweighted least squares

I) First Order update Rule

$$\theta_{t+1} = \theta_t - \alpha \nabla J(\theta)$$

Typically g : gradient $\nabla J(\theta)$

H : Hessian

Iteratively Reweighted least squares

II) Second Order update Rule

$$\theta_{t+1} = \theta_t - H^{-1} g$$

Iteratively Reweighted least squares

II) Second Order update Rule

$$\theta_{t+1} = \theta_t - H^{-1} g$$

For logistic regression:

$$\theta_{t+1} = \theta_t - (X^T \Delta X)^{-1} X^T (\hat{y} - y)$$

Iteratively Reweighted least squares

II) Second Order update Rule

$$\theta_{t+1} = \theta_t - H^{-1} g$$

For logistic regression:

$$\theta_{t+1} = \theta_t - (X^T D X)^{-1} X^T (\hat{y} - y)$$

$$= (X^T D X)^{-1} \left[X^T D X \theta_t - X^T \hat{y} + y \right]$$

$$\theta_{t+1} = (X^T D X)^{-1} X^T D \left[X \theta_t - \hat{y} + y \right]$$

Iteratively Reweighted least squares

II) Second Order update Rule

$$\begin{aligned}\theta_{t+1} &= (X^T D X)^{-1} X^T D [X \theta_t - \bar{D}^{-1} (\hat{y} - y)] \\ &= (X^T D X)^{-1} X^T D z_t \quad ; z_t = X \theta_t - \bar{D}^{-1} (\hat{y} - y)\end{aligned}$$

Contrast w/ weighted linear regression

$$\hat{\theta} = (X^T D X)^{-1} X^T D y$$

Iteratively Reweighted least squares

II) Second Order update Rule

$$\begin{aligned}\theta_{t+1} &= (X^T D X)^{-1} X^T D [X \theta_t - \bar{D}^{-1} (\hat{y} - y)] \\ &= (X^T D X)^{-1} X^T D z_t \quad ; z_t = X \theta_t - \bar{D}^{-1} (\hat{y} - y)\end{aligned}$$

Contrast w/ weighted linear regression

$$\hat{\theta} = (X^T D X)^{-1} X^T D y$$

Iteratively Reweighted least squares

II) Second Order update Rule

$$\begin{aligned}\theta_{t+1} &= (X^T D X)^{-1} X^T D [X \theta_t - \bar{\sigma}^{-1} (\hat{y} - y)] \\ &= (X^T D X)^{-1} X^T D z_t \quad ; z_t = X \theta_t - \bar{\sigma}^{-1} (\hat{y} - y)\end{aligned}$$

Contrast w/ weighted linear regression

$$\hat{\theta} = (X^T D X)^{-1} X^T D y$$

Iteratively Reweighted least squares

II) Second Order update Rule

$$\theta_{t+1} = (X^T D X)^{-1} X^T D [X \theta_t - \bar{\sigma}^{-1} (\hat{y} - y)]$$
$$= (X^T D X)^{-1} X^T D z_t \quad ; z_t = X \theta_t - \bar{\sigma}^{-1} (\hat{y} - y)$$

Contrast w/ weighted linear regression

$$\hat{\theta} = (X^T D X)^{-1} X^T D y$$

$$D = \begin{bmatrix} \bar{\sigma}_1^{-1} (\hat{y}_1 - y_1) \\ \vdots \\ \vdots \end{bmatrix}$$

Iteratively Reweighted least squares

II) Second Order update Rule

$$\theta_{t+1} = (X^T D X)^{-1} X^T D [X \theta_t - \bar{D}^{-1} (\hat{y} - y)]$$
$$= (X^T D X)^{-1} X^T D z_t \quad ; z_t = X \theta_t - \bar{D}^{-1} (\hat{y} - y)$$

Contrast w/ weighted linear regression

$$\hat{\theta} = (X^T D X)^{-1} X^T D y$$

$$D = \begin{bmatrix} \hat{y}_1 (1 + \hat{y}_1^2) & & \\ & \ddots & \\ & & \ddots \end{bmatrix}$$

$f(\theta_t)$ points to $\hat{y}_1 (1 + \hat{y}_1^2)$