

$$J : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$$

$$\theta \mapsto J(\theta)$$

Sea  $j = 0, 1, \dots, n$

$$\frac{\partial J}{\partial \theta_j} (\theta) = \left( \frac{1}{m} \sum_{i=1}^m \left[ h_\theta(x^{(i)}) - y^{(i)} \right] x_j^{(i)} \right)$$

$$\mu_i = h_\theta(x^{(i)}) = \sigma(\theta^T x^{(i)})$$

$x^{(i)}$  es la  $i$ -ésima instancia ( $\in \mathbb{R}^{n+1}$ )

$$\mu := \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_m \end{bmatrix} \quad y := \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\text{Diagram illustrating the cost function gradient calculation:}$$

$$\frac{\partial J}{\partial \theta_j} (\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\text{where } h_\theta(x) = \sigma(\theta^T x)$$

$$\theta = \begin{bmatrix} \alpha_1 & \alpha_2 & \alpha_3 \\ \beta_1 & \beta_2 & \beta_3 \end{bmatrix}$$

$$x = \begin{bmatrix} x_0^{(1)} & \dots & x_0^{(m)} \\ \vdots & \ddots & \vdots \\ x_n^{(1)} & \dots & x_n^{(m)} \end{bmatrix}_{(n+1) \times m}$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}_{m \times 1}$$

$$\frac{\partial J}{\partial \theta_j} (\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$g := \begin{bmatrix} \frac{\partial J}{\partial \theta_0} (\theta) \\ \vdots \\ \frac{\partial J}{\partial \theta_n} (\theta) \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m (\mu_i - y^{(i)}) x_0^{(i)} \\ \vdots \\ \sum_{i=1}^m (\mu_i - y^{(i)}) x_n^{(i)} \end{bmatrix}$$

$$= \sum_{i=1}^m (\mu_i - y^{(i)}) \begin{bmatrix} x_0^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$$

$$= X^T (\mu - y)$$

BLANK

Sean  $j, k = 0, 1, \dots, n$

$$\frac{\partial^2 J}{\partial \theta_k \partial \theta_j}(\theta) = \frac{\partial}{\partial \theta_k} \left( \frac{\partial J}{\partial \theta_j}(\theta) \right)$$

$$= \frac{\partial}{\partial \theta_k} \left( \sum_{i=1}^m (\mu_i - y_i) x_j^{(i)} \right)$$

$$= \sum_{i=1}^m \frac{\partial}{\partial \theta_k} x_j^{(i)} \sigma(\theta^\top x^{(i)})$$

$$= \sum_{i=1}^m x_j^{(i)} \sigma(\theta^\top x^{(i)}) \left[ 1 - \sigma(\theta^\top x^{(i)}) \right] x_k^{(i)}$$

$$= \sum_{i=1}^m x_j^{(i)} \mu_i (1 - \mu_i) x_k^{(i)}$$

$$= \sum_{i=1}^m \mu_i (1 - \mu_i) x_j^{(i)} x_k^{(i)}$$

?

$$\left[ \begin{array}{c} \frac{\partial^2 J}{\partial \theta_k \partial \theta_j} \end{array} \right]_{jk}$$

$$= \left[ \begin{array}{cccc} x_0^{(1)} & \cdots & x_0^{(m)} \\ \vdots & \ddots & \vdots \\ \mu_1 (1 - \mu_1) x_j^{(1)} & \cdots & \mu_m (1 - \mu_m) x_j^{(m)} \\ \vdots & \ddots & \vdots \\ x_n^{(1)} & \cdots & x_n^{(m)} \end{array} \right]_{(n+1) \times m} \left[ \begin{array}{c} x_0^{(1)} \\ \vdots \\ x_0^{(m)} \\ \vdots \\ x_k^{(1)} \\ \vdots \\ x_k^{(m)} \\ \vdots \\ x_n^{(1)} \\ \vdots \\ x_n^{(m)} \end{array} \right]_{m \times (n+1)}$$

$$= \left[ \begin{array}{cc} x_0^{(1)} & \cdots x_0^{(m)} \\ \vdots & \vdots \\ x_n^{(1)} & \cdots x_n^{(m)} \end{array} \right]_{(n+1) \times m} \left[ \begin{array}{c} \mu_1 (1 - \mu_1) \\ \vdots \\ \mu_m (1 - \mu_m) \end{array} \right]_{m \times 1}$$

$$= \cancel{X^T S X} \quad \vdash : \mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1} \times \mathbb{R}^{n+1}$$

$$x^T A x$$

$$\theta \mapsto \boxed{H(\theta)}_{n+1}$$

BLANK

$$J : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$$

$$H : \mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1} \times \mathbb{R}^{n+1}$$

$$H = X^T S X$$

donde

$$S = \begin{bmatrix} >0 \\ \mu_1(1-\mu_1) & \dots & & \\ & \dots & & \mu_m(1-\mu_m) \\ & & & >0 \end{bmatrix}$$

$$X \underset{M \times (n+1)}{\sim} \quad (M > n+1)$$

Dicir que  $X$  es de  $\textcircled{3}$   
rango máximo (full rank)

es decir que  $\text{rank}(X) = n+1$

$$\text{Sea } a = \begin{bmatrix} \alpha_0 \\ \vdots \\ \alpha_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

Como  $X$  es de rango máximo, si

$$\alpha_0 \begin{bmatrix} X^0 \\ \vdots \\ X^n \end{bmatrix} + \dots + \alpha_n \begin{bmatrix} X^0 \\ \vdots \\ X^n \end{bmatrix} = \begin{bmatrix} X^0 \dots X^n \end{bmatrix} \begin{bmatrix} \alpha_0 \\ \vdots \\ \alpha_n \end{bmatrix} = Xa = 0, \text{ if } b$$

entonces  $\alpha_0 = \alpha_1 = \dots = \alpha_n = 0$ , i.e.,  $a = 0$

Luego, esto es equivalente a  
 $\forall a \in \mathbb{R}^{n+1} \setminus \{0\} : Xa \neq 0$

Finalmente, veamos que  $H$  sea definida positiva.

En efecto, sea  $a \in \mathbb{R}^{n+1} \setminus \{0\}$

$$a^T H a = a^T X^T S X a = (Xa)^T S (Xa)$$

$$= \sum_{i=1}^m \underbrace{\mu_i(1-\mu_i)}_{>0} b_i^2 > 0$$

## Exercise 8.4

(a) For every  $k=1, 2, \dots, C$  let

$$v_k : \mathbb{R}^C \rightarrow \mathbb{R}$$

$$\begin{bmatrix} n_1 \\ \vdots \\ n_C \end{bmatrix} = \eta \mapsto v_k(\eta) := \frac{\exp(n_k)}{\sum_{c=1}^C \exp(n_c)}.$$

(4)

Let  $k, j = 1, 2, \dots, C$ . Show that

$$\frac{\partial v_k}{\partial n_j}(\eta) = v_k(\eta) (\delta_{kj} - v_j(\eta)).$$

Indeed, let  $\eta \in \mathbb{R}^C$ ; then

$$\begin{aligned} \frac{\partial v_k}{\partial n_j}(\eta) &= \frac{\frac{\partial}{\partial n_j} (\exp(n_k)) \sum - \exp(n_k) \sum_{c=1}^C \frac{\partial}{\partial n_j} (\exp(n_c))}{\sum^2} \\ &= \frac{\delta_{kj} \exp(n_k) \sum - \exp(n_k) \exp(n_j)}{\sum^2} \\ &= \frac{\exp(n_k)}{\sum} \left( \delta_{kj} - \frac{\exp(n_j)}{\sum} \right) \\ &= v_k(\eta) (\delta_{kj} - v_j(\eta)) \end{aligned}$$

(b) Let  $l: \mathbb{R}^{D \times C} \rightarrow \mathbb{R}$

$$[w_1 \dots w_C] = w \mapsto l(w) := \sum_{i=1}^m \sum_{k=1}^C y_{ik} \log(\mu_{ik}(w))$$

where  $y_{ik} := \delta_{y_{ik}}, \quad \mu_{ik}: \mathbb{R}^{D \times C} \rightarrow \mathbb{R}$

$$w \mapsto \frac{\exp(w_k^T x_i)}{\sum_{c'=1}^C \exp(w_{c'}^T x_i)}$$

and  $x_1, x_2, \dots, x_m$  are the instances.

Hence, show that

$$\nabla_w l = \begin{bmatrix} \nabla_{w_1} l \\ \vdots \\ \nabla_{w_C} l \end{bmatrix}$$

where

$$\nabla_{w_c} l (w) = \sum_{i=1}^m (y_{ic} - \mu_{ic}(w)) x_i$$

for every  $c = 1, 2, \dots, C$ .

5

Indeed, let  $c=1, 2, \dots, C$ . Then

6

$$\begin{aligned}
 \nabla_{w_c} l(w) &= \left[ \begin{array}{c} \frac{\partial}{\partial w_{c1}} (l(w)) \\ \vdots \\ \frac{\partial}{\partial w_{cD}} (l(w)) \end{array} \right] \\
 &= \left[ \begin{array}{c} \sum_{i=1}^m \sum_{k=1}^C y_{ik} \frac{1}{\mu_{ik}(w)} \frac{\partial \mu_{ik}}{\partial w_{c1}} (w) \\ \vdots \\ \sum_{i=1}^m \sum_{k=1}^C y_{ik} \frac{1}{\mu_{ik}(w)} \frac{\partial \mu_{ik}}{\partial w_{cD}} (w) \end{array} \right] \\
 &= \sum_{i=1}^m \sum_{k=1}^C \frac{y_{ik}}{\mu_{ik}(w)} \left[ \begin{array}{c} \frac{\partial \mu_{ik}}{\partial w_{c1}} (w) \\ \vdots \\ \frac{\partial \mu_{ik}}{\partial w_{cD}} (w) \end{array} \right]
 \end{aligned}$$

Let  $i = 1, 2, \dots, m$ . Then, let  $g_i : \mathbb{R}^{D \times G} \rightarrow \mathbb{R}^G$

$$w \mapsto g_i(w) := \begin{bmatrix} w_1^\top x_i \\ \vdots \\ w_G^\top x_i \end{bmatrix} =: \begin{bmatrix} \eta_1(w) \\ \vdots \\ \eta_G(w) \end{bmatrix}$$

Now, let  $k = 1, 2, \dots, G$  and let  $v_k$  as item (a) above.

Since the following diagram holds:

$$\begin{array}{ccc} \mathbb{R}^{D \times G} & \xrightarrow{g_i} & \mathbb{R}^G \\ & \searrow \mu_{ik} & \downarrow v_k \\ & & \mathbb{R} \end{array}$$

6  
7

We can apply 'the chain rule' once again! Then, for every  $r = 1, 2, \dots, D$

$$\begin{aligned} \frac{\partial v_k \circ g_i}{\partial w_{cr}}(w) &= \sum_{j=1}^G \frac{\partial v_k}{\partial \eta_j}(g_i(w)) \frac{\partial \eta_j}{\partial w_{cr}}(w) \\ &= \frac{\partial v_k}{\partial \eta_c}(g_i(w)) \frac{\partial \eta_c}{\partial w_{cr}}(w) \\ &\stackrel{\text{item(a)}}{=} v_k(g_i(w)) (\delta_{kc} - v_c(g_i(w))) x_{ir} \\ &= \mu_{ik}(w) (\delta_{kc} - \mu_{ic}(w)) x_{ir} \end{aligned}$$

CO

(8)

$$\nabla_{w_c} l(w) = \sum_{i=1}^m \sum_{k=1}^C \frac{y_{ik}}{\mu_{ik}(w)} \begin{bmatrix} \mu_{ik}(w)(\delta_{kc} - \mu_{ic}(w))x_{i1} \\ \mu_{ik}(w)(\delta_{kc} - \mu_{ic}(w))x_{iD} \end{bmatrix}$$

$$= \sum_{i=1}^m \sum_{k=1}^C y_{ik} (\delta_{kc} - \mu_{ic}(w)) \begin{bmatrix} x_{i1} \\ \vdots \\ x_{iD} \end{bmatrix}$$

$$= \sum_{i=1}^m \sum_{k=1}^C y_{ik} (\delta_{kc} - \mu_{ic}(w)) x_i$$

$$= \sum_{i=1}^m \left[ y_{ic} - y_{ic} \mu_{ic}(w) - \sum_{k \neq c} y_{ik} \mu_{ic}(w) \right] x_i$$

$$= \sum_{i=1}^m \left[ y_{ic} - \left( \sum_{k=1}^C y_{ik} \right) \mu_{ic}(w) \right] x_i$$

$$= \sum_{i=1}^m (y_{ic} - \mu_{ic}(w)) x_i$$

9

(c) Show that the Hessian of  $\ell : \mathcal{H}\ell$  is given by

$$\mathcal{H}\ell(w) = \begin{bmatrix} H_{1,1} \ell(w) & \cdots & H_{1,C} \ell(w) \\ \vdots & & \vdots \\ H_{C,1} \ell(w) & \cdots & H_{C,C} \ell(w) \end{bmatrix}$$

where

$$H_{c,c'} \ell(w) = - \sum_{i=1}^m \mu_{ic}(w) (\delta_{c,c'} - \mu_{i,c'}(w)) x_i x_i^T$$

for every  $c, c' = 1, 2, \dots, C$ .

Indeed, let  $c, c' = 1, 2, \dots, C$ . Then, by definition of the Hessian and item (b), the element  $(j, k)$  of the matrix  $H_{c,c'} \ell(w)$  is given by

$$\frac{\partial}{\partial w_{c',k}} \left[ \sum_{i=1}^m (y_{ic} - \mu_{ic}(w)) x_{i,j} \right] = - \sum_{i=1}^m x_{i,j} \frac{\partial v_c \circ g_i(w)}{\partial w_{c',k}}$$

Does NOT depend  
on  $w_{c',k}$

$\bullet \bullet$

$\mu_{ic} = v_c \circ g_i$   
 $g_i(w) = w^T x_i$

$$\begin{aligned}
 & - \sum_{i=1}^m x_{i,j} \left[ \sum_{s=1}^C \frac{\partial v_c}{\partial \eta_{i,s}} (g_i(w)) \frac{\partial \eta_{i,s}}{\partial w_{c',k}} \right] \\
 & = - \sum_{i=1}^m x_{i,j} \left[ \frac{\partial v_c}{\partial \eta_{i,c'}} (g_i(w)) \frac{\partial \eta_{i,c'}}{\partial w_{c',k}} \right]
 \end{aligned}$$

(9)

$$= - \sum_{i=1}^m x_{ij} \left[ \frac{\partial v_c}{\partial n_{i,c}} (g_i(w)) x_{ik} \right]$$

$$\stackrel{\text{by item (a)}}{=} - \sum_{i=1}^m x_{ij} v_c(g_i(w)) (\delta_{c,c} - v_{c'}(g_i(w))) x_{ik}$$

$$= - \sum_{i=1}^m \left[ \mu_{ic}(w) (\delta_{c,c} - \mu_{i,c'}(w)) \right] x_{ij} x_{ik}$$

Then, since  $\left[ \cdot \right]$  does NOT depend on  $j$  neither  $k$ ,

$$j, k \in \{1, 2, \dots, D\} \text{ and } x_i x_i^\top = \begin{bmatrix} x_{i1} x_{i1} \dots x_{i1} x_{iD} \\ \vdots \\ x_{iD} x_{i1} \dots x_{iD} x_{iD} \end{bmatrix},$$

$$H_{c,c'}(w) = - \sum_{i=1}^m \mu_{ic}(w) (\delta_{c,c} - \mu_{i,c'}(w)) x_i x_i^\top.$$

□

In  
the