

Problem 1:

a) For the multiclass logistic regression, we know,

$$P(\omega_k | \mathbf{x}) = y_k(\mathbf{x}) = \frac{\exp(a_k)}{\sum_i \exp(a_i)}$$

where,  $a_k = \omega_k' \mathbf{x}$

$$\begin{aligned} \text{So, } \frac{\partial y_k(\mathbf{x})}{\partial a_j} &= \frac{1}{\left(\sum_i e^{a_i}\right)^2} \left[ \sum_i e^{a_i} \frac{\partial e^{a_k}}{\partial a_j} - e^{a_k} \frac{\partial \sum_i e^{a_i}}{\partial a_j} \right] \\ &= \frac{1}{\left(\sum_i e^{a_i}\right)^2} \left[ e^{a_k} \sum_i e^{a_i} \frac{\partial a_k}{\partial a_j} - \frac{e^{a_k} e^{a_j}}{1} \right] \end{aligned}$$

let  $\delta_{k,j} = \begin{cases} 1 & \text{if } k=j \\ 0 & \text{o.w.} \end{cases}$

$$\begin{aligned} &= \frac{1}{\left(\sum_i e^{a_i}\right)^2} \left[ \delta_{k,j} e^{a_k} \sum_i e^{a_i} - e^{a_k} e^{a_j} \right] \\ &= \frac{e^{a_k}}{\sum_i e^{a_i}} \left[ \delta_{k,j} - \frac{e^{a_j}}{\sum_i e^{a_i}} \right] \\ &= y_k (\delta_{k,j} - y_j) \quad (\text{Proved}) \end{aligned}$$

(2)

Problem 1:

b)

we have,

$$P(\omega_k | x_i) = \gamma_k(x_i) = \frac{e^{a_{ik}}}{\sum_j e^{a_{ij}}}$$

Where,  $a_{ij} = \omega_j' x_i$  (activation function)

Also, we know, the cross-entropy error function

$$E(W) = - \sum_{i=1}^n \sum_{k=1}^c t_{ik} \ln \gamma_{ik}$$

Where,  $t_{ik} = 1$  if one-of-c encoding  $(x_i, \omega_k)$  is  $(x_i, t_i)$  and  $t_{ij} = 0$  if  $j \neq k$ .

$$\text{So, } \frac{\partial E(W)}{\partial \omega_j} = \sum_i \frac{\partial E(W)}{\partial a_{ij}} \frac{\partial a_{ij}}{\partial \omega_j} \quad (\text{By chain rule})$$

$$\text{Now, } \frac{\partial a_{ij}}{\partial \omega_j} = \frac{\partial}{\partial \omega_j} (\omega_j' x_i) = x_i \quad \text{--- (1)}$$

$$\text{and } \frac{\partial E(W)}{\partial a_{ij}} = - \frac{\partial}{\partial a_{ij}} \left[ \sum_{i=1}^n \sum_{k=1}^c t_{ik} \ln \gamma_{ik} \right]$$

$$= - \sum_{k=1}^c \frac{\partial}{\partial a_{ij}} (t_{ik} \ln \gamma_{ik})$$

$$= - \sum_{k=1}^c t_{ik} \frac{1}{\gamma_{ik}} \frac{\partial \gamma_{ik}}{\partial a_{ij}}$$

$$= - \sum_{k=1}^c t_{ik} \frac{1}{\gamma_{ik}} \cdot \gamma_{ik} (\delta_{kj} - \gamma_{ij}) \quad (\text{from part a})$$

$$= - \sum_{k=1}^c t_{ik} \delta_{kj} + \gamma_{ij} \sum_{k=1}^c t_{ik}$$

Since, if  $k=j$  then $\delta_{kj} = 1$  otherwise 0.Similarly for  $t_{ik}$ .

$$= \gamma_{ij} - t_{ij} \quad \text{--- (2)}$$

So, we can write gradient equation — form of (1) & (2)

$$\frac{\partial E(W)}{\partial \underline{w}_j} = \sum_{i=1}^n (Y_{ij} - t_{ij}) \underline{x}_i$$

Thus, Batch gradient descent rule can be written as —

update  $\underline{w}_j \leftarrow \underline{w}_j - \eta \sum_{i=1}^n (Y_{ij} - t_{ij}) \underline{x}_i$

and for each update, (single update)

$$\underline{w}_j \leftarrow \underline{w}_j - \eta (Y_{ij} - t_{ij}) \underline{x}_i$$