

IFT-6390 Fundamentals of Machine Learning

Professor: Ioannis Mitliagdas

Students: Abhay Luri (20209505), Saurabh Bodhe (20208545)

Homework 3 - Theoretical Part

Solutions:

Q) $(\text{ReLU}) g(x) = \max\{0, x\}$

this $\max\{0, x\}$ can be expanded as

$$\max(0, x) = \begin{cases} 0, & x < 0 \\ x, & x \geq 0 \end{cases}$$

Therefore, the first order derivative can be written as

$$\frac{d}{dx}(\text{ReLU})g(x) = \frac{d}{dx} \max(0, x) = \frac{d}{dx} \begin{cases} 0, & x < 0 \\ x, & x \geq 0 \end{cases}$$

and can be simplified as

$$\frac{d}{dx}(\text{ReLU})g(x) = \begin{cases} \frac{d}{dx} 0, & x < 0 \\ \frac{d}{dx} x, & x \geq 0 \end{cases}$$

Here individual derivatives can be solved as follows:

$$\frac{d}{dx} x^n = n x^{n-1}, \text{ and } \frac{d}{dx}(0) = \text{rect}'$$

$$\frac{d}{dx}(\text{ReLU})g(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

$$b) \sigma(x) = \frac{1}{1+e^{-x}}$$

$$= (1+e^{-x})^{-1}$$

$$\text{also } (1+e^{-x}) = \frac{1}{\sigma(x)} \quad \text{--- } \textcircled{1}$$

$$e^{-x} = \frac{1}{\sigma(x)} - 1$$

$$e^{-x} = \frac{1-\sigma(x)}{\sigma(x)} \quad \text{--- } \textcircled{2}$$

Therefore, the first order derivative

$$\sigma'(x) = -1 \times (1+e^{-x})^{-2} \times e^{-x} \times (-1)$$

$$\sigma'(x) = e^{-x} \times \frac{1}{(1+e^{-x})^2}$$

$$\sigma'(x) = \frac{e^{-x}}{(1+e^{-x})^2} \quad \text{--- } \textcircled{3}$$

Putting $\textcircled{1}$ & $\textcircled{2}$ in $\textcircled{3}$ we get

$$\sigma'(x) = \frac{1-\sigma(x)}{\sigma(x)} \times \sigma^2(x)$$

$$\boxed{\sigma'(x) = \sigma(x)(1-\sigma(x))}$$

$$c) \sigma(x) = \frac{1}{2} (\tanh(\frac{1}{2}x) + 1)$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\tanh\left(\frac{x}{2}\right) = \frac{e^{x/2} - e^{-x/2}}{e^{x/2} + e^{-x/2}}$$

$$\tanh\left(\frac{1}{2}x\right) + 1 = \frac{e^{\frac{1}{2}x} - e^{-\frac{1}{2}x}}{e^{\frac{1}{2}x} + e^{-\frac{1}{2}x}} + 1$$

$$= \frac{2e^{\frac{x}{2}}}{e^{\frac{x}{2}} + e^{-\frac{x}{2}}}$$

$$\frac{1}{2}(\tanh\left(\frac{x}{2}\right) + 1) = \frac{e^{\frac{x}{2}}}{e^{\frac{x}{2}} + e^{-\frac{x}{2}}}$$

$$= \frac{e^{\frac{x}{2}}}{e^{\frac{x}{2}} \left(1 + \frac{e^{-\frac{x}{2}}}{e^{\frac{x}{2}}}\right)}$$

$$\boxed{\frac{1}{2}(\tanh\left(\frac{x}{2}\right) + 1) = \frac{1}{1+e^{-x}} = \sigma(x)} \quad \text{proved}$$

d) $\ln \sigma(x) = -\text{softplus}(-x)$

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

Taking \ln on both sides, we get

$$\ln \sigma(x) = \ln\left(\frac{1}{1+e^{-x}}\right)$$

$$\ln \sigma(x) = \cancel{\ln(1)}^0 - \ln(1+e^{-x})$$

$$\ln \sigma(x) = -\ln(1+e^{-x}) = -\text{softplus}(-x)$$

Hence Proved

$$e) \text{Softplus}(x) - \text{Softplus}(-x) = x$$

As we know that

$$= \ln(1+e^x) - \ln(1+e^{-x})$$

$$= \ln\left(\frac{1+e^x}{1+e^{-x}}\right).$$

$$= \ln\left(\frac{1+e^x}{1+\frac{1}{e^x}}\right)$$

$$= \ln\left(\frac{e^x(1+e^{-x})}{1+e^x}\right)$$

$$= x$$

$$f) \text{Sign}(x) = 1_{x>0}(x) - 1_{x<0}(x)$$

$$g) \|x\|_2^2 = \sum_i x_i^2$$

$$\frac{\partial \|x\|_2^2}{\partial x_i} = 2x_i$$

$$\frac{\partial \|x\|_2^2}{\partial x} = \begin{bmatrix} 2x_1 \\ 2x_2 \\ \vdots \\ 2x_n \end{bmatrix} = \mathbf{2x}$$

$$h) \quad \|x_i\| = \sum_i |x_i|$$

$$\frac{\partial \|x\|_1}{\partial x} = \frac{\partial \sum_i |x_i|}{\partial x}$$

$$= \begin{bmatrix} \frac{\partial |x_1|}{\partial x_1} \\ \vdots \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{\partial |x_2|}{\partial x_2} \\ \vdots \\ 0 \end{bmatrix} + \dots$$

$$= \begin{bmatrix} \frac{x_1}{|x_1|} \\ 0 \\ \vdots \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{x_2}{|x_2|} \\ \vdots \\ 0 \end{bmatrix} + \dots$$

$$= \begin{bmatrix} \frac{x_1}{|x_1|} \\ \frac{x_2}{|x_2|} \\ \vdots \\ 0 \end{bmatrix} = \text{sign}(x)$$

$$i) \quad S(x)_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$S(c \cdot x_i) = \frac{e^{c \cdot x_i}}{\sum_j e^{c \cdot x_j}} \neq \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Let us consider two cases \therefore

for $c = 0$,

$$S(c \cdot x_i) = \frac{e^{c \cdot x_i}}{\sum_j e^{c \cdot x_j}} = \frac{1}{N} \quad N \Rightarrow \text{No. of elements in } X.$$

for $c \rightarrow \infty$

$$\lim_{c \rightarrow \infty} S(c \cdot x_i) = \frac{e^{c \cdot x_i}}{\sum_j e^{c \cdot x_j}} = 1 \text{ for maximum value of } x$$

Hence, Softmax is not invariant under Scalar multiplication.

j) As we know that,

$$S(x)_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$S(x+c)_i = \frac{e^{x_i+c}}{\sum_j e^{x_j+c}} = \frac{e^{x_i} \cdot e^c}{\sum_j e^{x_j} \cdot e^c}$$

$$= \frac{e^c \cdot e^{x_i}}{e^c \cdot \sum_j e^{x_j}} \boxed{= \frac{e^{x_i}}{\sum_j e^{x_j}}}$$

Hence, Softmax is translation invariant

k) $\frac{\partial S(x)_i}{\partial x_j} = S(x)_i \delta_{ij} - S(x)_i S(x)_j$

Consider 2 cases, $i=j$ and $i \neq j$

Case 1 : $i = j$

$$\frac{\partial S(x)_i}{\partial x_i} = \frac{\partial}{\partial x_i} \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Applying chain rule = $\frac{1}{\sum_j e^{x_j}} \frac{\partial e^{x_i}}{\partial x_i} + e^{x_i} \frac{\partial}{\partial x_i} \frac{1}{\sum_j e^{x_j}}$

$$= \frac{e^{x_i}}{\sum_j e^{x_j}} - \frac{e^{x_i} e^{x_i}}{(\sum_j e^{x_j})^2}$$

$$= S(x)_i - S(x)_i^2$$

$$= S(x)_i - S(x)_i S(x)_i \quad \text{--- } ①$$

Case 2 : $i \neq j$

$$\frac{\partial S(x)_i}{\partial x_i} = \frac{\partial}{\partial x_i} \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$= \frac{1}{\sum_j e^{x_j}} \frac{\partial e^{x_i}}{\partial x_j} + e^{x_i} \frac{\partial}{\partial x_j} \frac{1}{\sum_j e^{x_j}}$$

$$= 0 - \frac{e^{x_i} e^{x_j}}{(\sum_j e^{x_j})^2}$$

$$= 0 - S(x)_i S(x)_j \quad \text{--- } ②$$

Combining ① & ② we get

$$\boxed{\frac{\partial S(x)_i}{\partial x_j} = S(x)_i \delta_{ij} - S(x)_i S(x)_j}$$

L) As given in the question

$$\frac{\partial S(x)_i}{\partial x_j} = S(x)_i \cdot 1_{i=j} - S(x)_i \cdot S(x)_j$$

With condition ($i=j$) diagonal elements of $\frac{\partial S(x)}{\partial x}$ are

$$\left(\frac{\partial S(x)}{\partial x} \right)_{i=j} = S(x)_i - S(x)_i \cdot S(x)_j$$

Otherwise $i \neq j$, we have

$$\left(\frac{\partial S(x)}{\partial x} \right)_{i \neq j} = -S(x)_i \cdot S(x)_j$$

We can write

$$\frac{\partial S(x)}{\partial x} = \text{diag}(S(x)) - \begin{bmatrix} S(x)_1 \\ S(x)_2 \\ \vdots \\ S(x)_n \end{bmatrix} \begin{bmatrix} S(x)_1 & S(x)_2 & \dots & S(x)_n \end{bmatrix}$$

$$\boxed{\frac{\partial S(x)}{\partial x} = \text{diag}(S(x)) - S(x)S(x)^T}$$

$$m) y = \sigma(x) = \begin{bmatrix} \sigma(x_1) \\ \sigma(x_2) \\ \vdots \\ \sigma(x_n) \end{bmatrix}$$

Jacobian of $f(x) = \sigma(x)$ is diagonal and can be written as

$$\left(\frac{\partial y}{\partial x} \right) = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \dots & \frac{\partial y_n}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \sigma(x_1)(1-\sigma(x_1)) \\ \vdots \\ \vdots \\ \sigma(x_n)(1-\sigma(x_n)) \end{bmatrix}$$

$$\nabla_x L = \left(\frac{\partial y}{\partial x} \right)^T \nabla_y L$$

$$\nabla_x L = \left(\frac{\partial y}{\partial x} \right)^T \nabla_y L = \underbrace{\begin{bmatrix} \sigma(x_1)(1-\sigma(x_1)) \\ \vdots \\ \vdots \\ \sigma(x_n)(1-\sigma(x_n)) \end{bmatrix}}_{d_{h2} \times d_{h2}} \underbrace{\begin{bmatrix} (\nabla_y L)_1 \\ (\nabla_y L)_2 \\ \vdots \\ (\nabla_y L)_n \end{bmatrix}}_{d_{h1} \times 1}$$

$$= \begin{bmatrix} \sigma(x_1)(1-\sigma(x_1)) (\nabla_y L)_1 \\ \vdots \\ \vdots \\ \sigma(x_n)(1-\sigma(x_n)) (\nabla_y L)_n \end{bmatrix} d_{h1} \times 1$$

for the case of $y = s(x)$, we have

$$\begin{aligned} \nabla_x L &= \left(\frac{\partial y}{\partial x} \right)^T \nabla_y L = \left(\text{diag}(s(x)) - s(x)s(x)^T \right)^T \nabla_y L \\ &= \left(\text{diag}(s(x)) - s(x)s(x)^T \right) \nabla_y L \\ &= \text{diag}(s(x)) \nabla_y L - s(x)s(x)^T \nabla_y L \\ &= \text{diag}(s(x)) \nabla_y L - s(x)(s(x)^T \nabla_y L) \end{aligned}$$

$$\underbrace{s(x)^T}_{1 \times n} \underbrace{\Delta y L}_{n \times 1} = O(n)$$

$$\underbrace{s(x)}_{n \times 1} \left(\underbrace{s(x)^T \nabla_y L}_{1 \times 1} \right) = O(n)$$

$$\text{diag}(s(x)) \nabla_y L = O(n)$$

The overall complexity is $O(n)$.

Q3.a) First Layer
Size = 32×32 , Num of kernels = 6, filter size = 5
Padding = 0, stride = 1
Output Dimension = $((32 - 5 + 2(0)) / 1) + 1$
 $= 27 + 1 = 28 \times 28 \times 6$

Second Layer

Size = 28×28 , pooling layer on this output : pooling layer dimension : 2×2 with no overlapping., stride = 2

$$\text{Output Dimension} = ((28 - 2 + 2(0)) / 2) + 1$$
$$= 14 \times 14 \times 6$$

Third Layer

Size = 14×14 , Num of kernels = 16, stride = 1,
padding = 0, filter size = 5

$$\text{Output Dimension} = ((14 - 5 + 2(0)) / 1) + 1$$
$$= 10 \times 10 \times 16$$

Dimensionality (scalar) output of last layer: $10 \times 10 \times 16$
 $= 1600$

b) parameters for last layer

$$= 5 \times 5 \times 16 \times 6$$

$$= 2400$$

c) input matrix = $64 \times 64 \times 3$

Output shape of last layer = $(64, 6, 6)$

Output Dimension = $\left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1$

in case of dilations,

$$\text{effective } k' = k + (k-1)(d-1)$$

Assume we are not using padding or dilation.

output shape = 6

Stride (s) = 1

Padding (p) = 0

Dilation (d) = 1

$$b = ((64 - k + 2(0)) / 1) + 1$$

$$b = 64 - k + 0 + 1$$

$$\begin{aligned} k &= 64 - b + 0 + 1 \\ \boxed{k} &= 59 \end{aligned}$$

d) Assume $d=2, p=1$

Output Shape = 6

Stride (s) = 1

$$b = ((64 - (2k-1) + 2(1)) / 1) + 1$$

$$b = ((64 - 2k + 1 + 2)) + 1$$

$$\boxed{k = 31}$$

e) Assume $d=1, p=1$

Output Shape = 6

Stride (s) = 1

$$b = ((64 - k + 2(1)) / 1) + 1$$

$$b = 64 - k + 2 + 1$$

$$\begin{aligned} k &= 64 - b + 2 + 1 \\ \boxed{k} &= 61 \end{aligned}$$

Qd a) Since $b^{(1)}$ contains one bias term for every neuron in the layer containing d_h neurons, its dimension is $d_h \times 1$.

$$h^a = W^{(1)}x + b^{(1)}$$

Here each row of $W^{(1)}$ and $b^{(1)}$ represents the weight and bias for a neuron.

$$\text{Output of the } j^{\text{th}} \text{ neuron} = h_j^a = W_j^{(1)}x + b_j^{(1)}$$

$$h^s = \phi(h^a) = \sin(h^a)$$

b) $W^{(2)}$ has dimensions (m, d_h) and $b^{(2)}$ dimension m .

Activation function of the neurons of the output layer

$$O^a = W^{(2)}h^s + b^{(2)}$$

$$O_k^a = W_k^{(2)}h^s + b_k^{(2)}$$

c)

$$O^s = \text{softmax}(O^a)$$

$$O_k^s = \text{softmax}(O_k^a) = \frac{e^{O_k^a}}{\sum_{j=1}^m e^{O_j^a}}$$

Since $e^x > 0$, it is evident that $O_k^s > 0$

$$\begin{aligned} \text{So, } \sum_{k=1}^m O_k^s &= \sum_{k=1}^m \text{softmax}(O_k^a) = \sum_{k=1}^m \frac{e^{O_k^a}}{\sum_{j=1}^m e^{O_j^a}} \\ &= \frac{e^{O_1^a} + e^{O_2^a} + \dots + e^{O_m^a}}{e^{O_1^a} + e^{O_2^a} + \dots + e^{O_m^a}} = 1 \end{aligned}$$

$$\begin{aligned}
 d) \quad L(x, y) &= -\log O_y^s(x) \\
 &= -\log \text{softmax}(O_y^a) \\
 &= -\log \frac{e^{O_y^a(x)}}{\sum_{j=1}^m e^{O_j^a(x)}}
 \end{aligned}$$

e) $\hat{R} \Rightarrow$ it is an estimate of the true error and is defined as the average loss over the dataset.

$$\begin{aligned}
 \hat{R} &= \frac{1}{n} \sum_{i=1}^N L(x^{(i)}, y^{(i)}) \\
 &= \frac{1}{n} \sum_{i=1}^N -\log \frac{e^{O_{y^{(i)}}^a(x^{(i)})}}{\sum_{j=1}^m e^{O_j^a(x^{(i)})}}
 \end{aligned}$$

Parameters of Θ are: $\{w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}\}$

Scalar parameters: n_Θ

$$\begin{aligned}
 n_\Theta &= d \cdot d_h + d_h + m \cdot d_h + m \\
 &= d_h(d+1) + m \cdot (d_h+1)
 \end{aligned}$$

The Optimization problem is .

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \hat{R}_\Theta$$

or

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^N L_\Theta(x^{(i)}, y^{(i)})$$

f) $\theta \leftarrow \theta - \eta \frac{d\hat{R}}{d\theta}$ $\eta \Rightarrow \text{Learning Rate}$
 $\theta \Rightarrow [w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}]$

for each iteration the update rule is shown below:

$$w^{(1)} := w^{(1)} - \eta \frac{d\hat{R}}{dw^{(1)}}$$

$$b^{(1)} := b^{(1)} - \eta \frac{d\hat{R}}{db^{(1)}}$$

$$w^{(2)} := w^{(2)} - \eta \frac{d\hat{R}}{dw^{(2)}}$$

$$b^{(2)} := b^{(2)} - \eta \frac{d\hat{R}}{db^{(2)}}$$

g) $\frac{\partial L}{\partial o^s} = -\frac{\partial \log o_y^s}{\partial o^s} = -o^{-s} \text{onehot}_m(y)$

$$= \begin{bmatrix} -\frac{1}{o_1^s} \\ \vdots \\ -\frac{1}{o_y^s} \\ \vdots \\ -\frac{1}{o_m^s} \end{bmatrix} - \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ -\frac{1}{o_y^s} \\ \vdots \\ 0 \end{bmatrix}$$

$$\frac{\partial o^s}{\partial o^a} = \frac{\partial \text{softmax}(o^a)}{\partial o^a}$$

$$\frac{\partial o^s}{\partial o^a} = \begin{bmatrix} o_1^s = o_1^s o_1^s & - & - & - & - & - & o_1^s o_m^s \\ | & & & & & & | \\ -o_m^s o_1^s & - & - & - & - & - & o_m^s o_m^s \end{bmatrix}$$

$$\frac{\partial L}{\partial o^a} = \frac{\partial L}{\partial o^s} \frac{\partial o^s}{\partial o^a}$$

$$= \begin{bmatrix} o_1^s - o_1^s o_1^s & \dots & \dots & o_1^s o_m^s \\ \vdots & & & \vdots \\ \vdots & & & \vdots \\ -o_m^s o_1^s & \dots & \dots & o_m^s - o_m^s o_m^s \end{bmatrix}_{m \times m} \begin{bmatrix} 0 \\ \vdots \\ -\frac{1}{o_y^s} \\ \vdots \\ 0 \end{bmatrix}_{m \times 1} = \begin{bmatrix} o_1^s \\ \vdots \\ o_y^s - 1 \\ \vdots \\ o_m^s \end{bmatrix}_{m \times 1}$$

$$= o^s - \text{Onehot}_m(y)$$

h) $o_k^a = w_k^{(2)} h^s + b_k^{(2)}$

$$\frac{\partial o_k^a}{\partial w_{kj}^{(2)}} = \frac{\partial w_k^{(2)} h^s}{\partial w_{kj}^{(2)}} + \frac{\partial b_k^{(2)}}{\partial w_j^{(2)}} = h_j^s \quad \rightarrow \textcircled{1}$$

$$\frac{\partial L}{\partial w_{kj}^{(2)}} = \frac{\partial L}{\partial o_k^a} \frac{\partial o_k^a}{\partial w_{kj}^{(2)}} \quad \rightarrow \textcircled{2}$$

Putting $\textcircled{1}$ in $\textcircled{2}$ we get

$$\boxed{\frac{\partial L}{\partial w_{kj}^{(2)}} = \frac{\partial L}{\partial o_k^a} h_j^s}$$

$$\frac{\partial o_k^a}{\partial b_k^{(2)}} = \frac{\partial w_k^{(2)} h^s}{\partial b_k^{(2)}} + \frac{\partial b_k^{(2)}}{\partial b_k^{(2)}} = 1 \quad \rightarrow \textcircled{3}$$

$$\frac{\partial L}{\partial b_k^{(2)}} = \frac{\partial L}{\partial o_k^a} \frac{\partial o_k^a}{\partial b_k^{(2)}} \quad \rightarrow \textcircled{4}$$

Putting $\textcircled{3}$ in $\textcircled{4}$ we get

$$\boxed{\frac{\partial L}{\partial b_k^{(2)}} = \frac{\partial L}{\partial o_k^a} \times (\textcircled{1})}$$

$$(i) \frac{\partial L}{\partial W^{(2)}} = \frac{\partial L}{\partial o^a} h^s T$$

$$\frac{\partial L}{\partial b^{(2)}} = \frac{\partial L}{\partial o^a}$$

Dimension of h^s : $d_h \times 1$

Dimension of $\frac{\partial L}{\partial W^{(2)}}$: $m \times d_h$

Dimension of $\frac{\partial L}{\partial b^{(2)}}$: $m \times 1$

$$j) \frac{\partial L}{\partial h_j^s} = \sum_{k=1}^m \frac{\partial L}{\partial o_k^a} \frac{\partial o_k^a}{\partial h_j^s}$$

$$o_k^a = W_k^{(2)} h^s + b_k^{(2)}$$

$$\frac{\partial o_k^a}{\partial h_j^s} = \frac{\partial W_k^{(2)} h^s}{\partial h_j^s} + \frac{\partial b_k^{(2)}}{\partial h_j^s} = W_{kj}^{(2)}$$

$$= \sum_{k=1}^m \frac{\partial L}{\partial o_k^a} \frac{\partial o_k^a}{\partial h_j^s} = \sum_{k=1}^m \frac{\partial L}{\partial o_k^a} W_{kj}^{(2)}$$

$$k) \frac{\partial L}{\partial h^s} = W^{(2)T} \frac{\partial L}{\partial o^a}$$

Dimension of h^s : $d_h \times 1$

Dimension of $\frac{\partial L}{\partial o^a}$: $m \times 1$

Dimension of $W^{(2)}$: $m \times d_h$

$$l) \frac{\partial L}{\partial h_j^a} = \frac{\partial L}{\partial h_j^s} \cdot \frac{\partial h_j^s}{\partial h_j^a}$$

$$h_j^s = \phi(h_j^a) = \sin(h_j^a)$$

$$\frac{\partial h_j^s}{\partial h_j^a} = \cos(h_j^a)$$

$$\frac{\partial L}{\partial h_j^a} = \frac{\partial L}{\partial h_j^s} \cos(h_j^a)$$

$$m) \frac{\partial L}{\partial h_j^a} = \frac{\partial L}{\partial h_j^s} \cos(h_j^a)$$

Dimension of $h_j^s = (d_h, 1)$

Dimension of $\frac{\partial L}{\partial h_j^s} = (d_h, 1)$

Dimension of $\frac{\partial L}{\partial h_j^a} = (d_h, 1)$

$$n) h_j^a = b_j^{(1)} + \sum_{i=1}^d w_{ji}^{(1)} x_i$$

chain rules gives:

$$\frac{\partial L}{\partial w_{kj}^{(1)}} = \frac{\partial L}{\partial h_k^a} \frac{\partial h_k^a}{\partial w_{kj}^{(1)}} = x_j$$

$$\frac{\partial L}{\partial b_j^{(1)}} = \frac{\partial L}{\partial h_j^a} \frac{\partial h_j^a}{\partial b_j^{(1)}} = 1$$

we have,

$$\frac{\partial L}{\partial w_{kj}^{(1)}} = \frac{\partial L}{\partial h_j^a} x_j$$

$$\frac{\partial L}{\partial b_j^{(1)}} = \frac{\partial L}{\partial h_j^a}$$

Q) $\frac{\partial L}{\partial w^{(1)}} = \frac{\partial L}{\partial h^a} x^T$

$$\frac{\partial L}{\partial w^{(1)}} = \frac{\partial L}{\partial h^a}$$

$\frac{\partial L}{\partial h^a}$ has dimension $(d_n, 1)$.

$\frac{\partial L}{\partial w^{(1)}}$ has dimension (d_n, d) .

x has dimension $(d, 1)$.

$$\frac{\partial L}{\partial b^{(1)}} = (d_n, 1)$$

P) We calculate $\frac{\partial L}{\partial h_j^a}$ in the previous question

and $h_j^a = b_j^{(1)} + \sum_{i=1}^d w_{j(i)}^{(1)} x_i$

chain rule gives:

$$\frac{\partial L}{\partial x_k} = \sum_j \frac{\partial L}{\partial h_j^a} \frac{\partial h_j^a}{\partial x_k}$$

$$\frac{\partial h_j^a}{\partial x_k} = w_{jk}^{(1)}$$

$$\frac{\partial L}{\partial x_k} = \sum_j \frac{\partial L}{\partial h_j^a} w_{jk}^{(1)}$$