

IFT6390 HW3 Theory Q1

| a) $\text{sign}(x) = \mathbb{1}_{x>0} - \mathbb{1}_{x<0}$

| b) $g(x) = \begin{cases} x, & \text{when } x > 0 \\ 0, & \text{when } x \leq 0 \end{cases}$

$$\nabla g(x) = \frac{d}{dx} g(x) = \begin{cases} 1, & \text{when } x > 0 \\ 0, & \text{when } x \leq 0 \end{cases}$$

Since the Heaviside step function is defined as:

$$H(x) = \begin{cases} 1, & \text{when } x > 0 \\ 0.5, & \text{when } x = 0 \\ 0, & \text{when } x < 0 \end{cases}$$

then $\nabla g(x) = H(x)$ for $x < 0$ and $x > 0$.

| c) $g(x) = x H(x)$

Alternatively.

$$g(x) = (x+1)^{H(x)} - 1$$

| d) $a(x) = \frac{1}{1 + e^{-kx}}$

when k is large.

$$\lim_{k \rightarrow \infty} a(x) = \lim_{k \rightarrow \infty} \left(\frac{1}{1 + e^{-kx}} \right)$$

when $x > 0$,

$$\begin{aligned} \lim_{k \rightarrow \infty} \left(\frac{1}{1 + e^{-kx}} \right) &= \lim_{k \rightarrow \infty} \left(\frac{1}{1+0} \right) \\ &= 1 \end{aligned}$$

when $x = 0$,

$$\begin{aligned} \lim_{k \rightarrow \infty} \left(\frac{1}{1 + e^{-kx}} \right) &= \lim_{k \rightarrow \infty} \left(\frac{1}{1+1} \right) \\ &= \frac{1}{2} \end{aligned}$$

when $x < 0$,

$$\begin{aligned} \lim_{k \rightarrow \infty} \left(\frac{1}{1 + e^{-kx}} \right) &= \lim_{k \rightarrow \infty} \left(\frac{1}{1+\infty} \right) \\ &= 0 \end{aligned}$$

Since at large k ,

$$a(x) = \begin{cases} 1, & \text{when } x > 0 \\ 0.5, & \text{when } x = 0 \\ 0, & \text{when } x < 0 \end{cases}$$

then $H(x)$ can be approximated by $a(x)$
when k is large.

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

$$\begin{aligned}
&= -\left(1+e^{-x}\right)^{-2} \left(-e^{-x}\right) \\
&= \frac{1}{1+e^{-x}} \left(\frac{e^{-x}}{1+e^{-x}}\right) \\
&= \frac{1}{1+e^{-x}} \left(\frac{1+e^{-x}-1}{1+e^{-x}}\right) \\
&= \frac{1}{1+e^{-x}} \left(1 - \frac{1}{1+e^{-x}}\right) \\
&= \sigma(x) \left(1 - \sigma(x)\right)
\end{aligned}$$

Let $x \in \mathbb{R}^n$, hence $\sigma(x) : \mathbb{R}^n \rightarrow \mathbb{R}$

$$J_{\sigma(x)} = \left[\frac{\partial \sigma}{\partial x_1}, \frac{\partial \sigma}{\partial x_2}, \dots, \frac{\partial \sigma}{\partial x_n} \right]$$

$$= \begin{bmatrix} \frac{\partial \sigma(x_1)}{\partial x_1} & \frac{\partial \sigma(x_1)}{\partial x_2} & \dots & \frac{\partial \sigma(x_1)}{\partial x_n} \\ \frac{\partial \sigma(x_2)}{\partial x_1} & \frac{\partial \sigma(x_2)}{\partial x_2} & \dots & \frac{\partial \sigma(x_2)}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \sigma(x_n)}{\partial x_1} & \frac{\partial \sigma(x_n)}{\partial x_2} & \dots & \frac{\partial \sigma(x_n)}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$= \begin{bmatrix} \sigma(x_1)(1-\sigma(x_1)) & 0 & \dots & 0 \\ 0 & \sigma(x_2)(1-\sigma(x_2)) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma(x_n)(1-\sigma(x_n)) \end{bmatrix}$$

$$\Rightarrow (J_{\sigma(x)})_{i,j} = \sigma(x_i)(1-\sigma(x_i))\delta_{i=j}, \text{ where}$$

$$i \in \{1, 2, \dots, n\}$$

$$j \in \{1, 2, \dots, n\}$$

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

$$\text{LHS} = \ln \left(\frac{1}{1 + \exp(-x)} \right) \Rightarrow -\ln(1 + \exp(-x)) \\ \Rightarrow -\ln(\text{softplus}(-x)) = \text{RHS}$$

$$g) \text{softplus}(x) - \text{softplus}(-x) = x$$

$$\begin{aligned} \text{LHS} &= \ln(1 + \exp(x)) - \ln(1 + \exp(-x)) \\ &= \ln(1 + \exp(x)) - \ln\left(1 + \frac{1}{\exp(x)}\right) \\ &= \ln\left(1 + \exp(x)\right) - \ln\left(\frac{1 + \exp(x)}{\exp(x)}\right) \\ &= \ln\left(\frac{(1 + \exp(x)) \times \exp(x)}{(1 + \exp(x))}\right) \\ &= \ln(\exp(x)) \\ &= x = \text{RHS} \end{aligned}$$

$$h) S(x+c) = S(x)$$

$$\text{LHS} = S(x+c)_i = \frac{\exp(x_i + c)}{\sum_j \exp(x_j + c)} = \frac{\exp(x_i) \exp(c)}{\cancel{\exp(c)} \sum_j \exp(x_j)} = \frac{\exp(x_i)}{\sum_j \exp(x_j)} = S(x) = \text{RHS}$$

$$i) S(x^c)$$

$$\text{LHS} = S(x^c)_i = \frac{\exp(x_i^c)}{\sum_j \exp(x_j^c)} = \frac{\exp(x_i)^c}{\sum_j \exp(x_j)^c}$$

$$\text{if } c=0 \Rightarrow S(x^c) = \frac{1}{\sum 1} = \frac{1}{K}, \text{ where } K = \text{no. of classes}$$

$$c \rightarrow 1 \Rightarrow S(x^c) = S(x)$$

As c increases, $S(x^c)$ gets more skewed

j) $\frac{\partial s(x)_i}{\partial x_j} = s(x_i) (s_{i=j} - s(x)_j)$

 $s(x_i) = \frac{\exp(x_i)}{\sum \exp(x_j)} \Rightarrow \partial s(x_j) = \exp(x_i) \times \frac{-1}{(\sum \exp(x_j))^2} \times (\exp(x_j))$
 $= \frac{\exp(x_i)}{\sum \exp(x_j)} \times -\frac{\exp(x_j)}{\sum \exp(x_j)}$
 $= s(x)_i \times -s(x)_j \text{ if } j \neq i$
 $\text{or } = s(x)_i (1 - s(x)_i) \text{ if } j = i$
 $= s(x)_i (s_{i=j} - s(x)_j) \text{ - RHS}$

k) Let Jacobian of softmax = J

$\Rightarrow J_{ij} = \frac{\partial s(x)_i}{\partial x_j} \text{ if } i=j \Rightarrow J_{ii} = \frac{\partial s(x)_i}{\partial x_i} = s(x)_i (-s(x)_i) = s(x) (-s(x))^\top$
 $i \neq j \Rightarrow J_{ij} = \frac{\partial s(x)_i}{\partial x_j} = s(x)_i (1 - s(x)_j) = s(x)_i (1 - s(x))^\top$
 $\therefore J = s(x) (1 - s(x))^\top - \text{diag}(s(x))$

l) $\nabla_u \log s(x(u)) = \nabla_u x(u)_i - E_j [\nabla_u x(u)_j]$

$\log s(x(u))_i = \log(\exp(x(u)_i)) - \log(\sum_j \exp(x(u)_j))$
 $= x(u)_i - \log(\sum_j \exp(x(u)_j))$

$\nabla_u \log s(x(u))_i = \nabla_u x(u)_i - \nabla_u \log(\sum_j \exp(x(u)_j))$
 $= \nabla_u x(u)_i - \frac{1}{\sum_j \exp(x(u)_j)} \times \sum_j \exp(x(u)_j) \nabla_u x(u)_j \quad \text{--- (1)}$

$E_j [\nabla_u x(u)_j] = \sum_i \nabla_u x(u)_j s(x(u))_j = \sum_i \nabla_u x(u)_j \frac{\exp(x(u)_j)}{\sum \exp(x(u)_j)} \quad \text{--- (2)}$

$\Rightarrow \nabla_u \log s(x(u))_i = \nabla_u x(u)_i - \frac{\sum_i \nabla_u x(u)_j \exp(x(u)_j)}{\sum \exp(x(u)_j)}$
 $= \nabla_u x(u)_i - E_j [\nabla_u x(u)_j]$

$$m) \quad L(x, c) = \sum_{i=1}^k -c_i \log y_i \quad \text{where } y = S(x)$$

Gradient of cross-entropy loss,

$$\begin{aligned}\nabla_u L(x, c) &= \sum_{i=1}^k -c_i \nabla_u \log S(x(u))_i \\ &= \sum_{i=1}^k -c_i (\nabla_u x(u)_i - E_j [\nabla_u x(u)_j]) \\ &= -\sum_{i=1}^k c_i \nabla_u x(u)_i + E_j [\nabla_u x(u)_j] \sum_{i=1}^k c_i\end{aligned}$$

Since c is a one-hot vector, $\sum c_i = 1$

$$\Rightarrow \nabla_u L(x, c) = \sum c_i \nabla_u x(u)_i + E_j [\nabla_u x(u)_j] \times 1$$

Q2

a) $W^{(1)} \in \mathbb{R}^{d_h \times d}$

Dimension of $b^{(1)} = d_h \times 1$

Formula of $h^a = W^{(1)}x + b^{(1)}$
where $x \in \mathbb{R}^d$

$$h_j^a = \sum_{k=1}^d w_{jk}^{(1)} x_k + b_j^{(1)}$$

$$h^s = \text{softplus}(h^a) \quad \text{where} \quad \begin{cases} h^s \in \mathbb{R}^{d_h} \\ h^a \in \mathbb{R}^{d_h} \end{cases}$$
$$h_j^s = \ln(1 + \exp(h_j^a))$$

$$h^s = \left[\ln(1 + \exp(h_1^a)) \quad \ln(1 + \exp(h_2^a)) \quad \dots \quad \ln(1 + \exp(h_{d_h}^a)) \right]^T$$

$$b) \text{ Dimensions of } W^{(2)} = m \times d_h$$

$$b^{(2)} = m \times 1$$

The pre-activation output vector denoted by O^a is defined as:

$$O^A = W^{(2)} h^s + b^{(2)} ; O^A \in \mathbb{R}^m$$

$$\text{Activation function} = \phi(O^A) = \text{softmax}(O^A)$$

$$O^s = \phi(O^A) = \text{softmax}(W^{(2)} h^s + b^{(2)})$$

$$\text{where, } O^s \in \mathbb{R}^m$$

$$O_k^a = \sum_{j=1}^{d_h} w_{kj}^{(2)} h_j^s + b_k^{(2)}$$

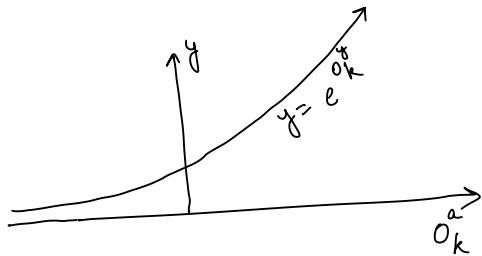
$$O_k^s = \frac{e^{O_k^a}}{\sum_{l=1}^m e^{O_l^a}}$$

$$c) \quad o^s = \text{softmax}(o^a) \quad ; \quad o^s \in \mathbb{R}^m$$

$$o_k^s = \frac{e^{o_k^a}}{\sum_{j=1}^m e^{o_j^a}} \quad \forall k \in \{1, 2 \dots m\}$$

$$\text{if } o_k^a \in \mathbb{R}, \quad e^{o_k^a} > 0$$

$$\text{therefore, } \sum_{j=1}^m e^{o_j^a} > 0$$



Since both numerator & denominator term are > 0

$$\text{hence } \forall k \in \{1, 2 \dots m\} \quad o_k^s > 0$$

$$\begin{aligned} \sum_{k=1}^m o_k^s &= \sum_{k=1}^m \left(\frac{e^{o_k^a}}{\sum_{j=1}^m e^{o_j^a}} \right) \\ &= \frac{1}{\sum_{j=1}^m e^{o_j^a}} \left(\sum_{k=1}^m e^{o_k^a} \right) = \frac{\sum_{k=1}^m e^{o_k^a}}{\sum_{j=1}^m e^{o_j^a}} = 1 \end{aligned}$$

Hence the sum of the vector $o^s = 1$

Since each element o_k^s represents probability values, their sum must be equal to 1 and should be individually positive.

d) A binary classifier with softmax activation would consist of 2 output neurons such that,

$$O_1^S = P(y=0) = \frac{e^{O_1^A}}{e^{O_1^A} + e^{O_2^A}} = \frac{1}{1 + e^{(O_2^A - O_1^A)}}$$

$$O_2^S = P(y=1) = \frac{e^{O_2^A}}{e^{O_1^A} + e^{O_2^A}} = \frac{1}{1 + e^{(O_1^A - O_2^A)}}$$

Both probability values are of the form

$$O_k^S = \frac{1}{1 + e^{-z}}$$

which is same as sigmoid activation function.

$$\text{and, } O_2^S = 1 - O_1^S$$

Thus it can be concluded :

- softmax activation is same as sigmoid for binary classification.
- softmax is an extension of sigmoid for multi class classification problems.

$$e) L_{\text{mse}}(\sigma(o^a), y) = (\sigma(o^a) - y)^2$$

where $L_{\text{mse}} \in \mathbb{R}$, $y \in \{0, 1\}$

For binary classification, $o^a \in \mathbb{R}$

$$\text{therefore } \frac{\partial L_{\text{mse}}}{\partial o^a} \in \mathbb{R}$$

$$\frac{\partial L_{\text{mse}}}{\partial o^a} = \frac{\partial (\sigma(o^a) - y)^2}{\partial o^a}$$

$$= 2(\sigma(o^a) - y) \frac{\partial \sigma(o^a)}{\partial o^a}$$

where,

$$\frac{\partial \sigma(o^a)}{\partial o^a} = \sigma(o^a)(1 - \sigma(o^a))$$

thus,

$$\boxed{\frac{\partial L_{\text{mse}}}{\partial o^a} = 2(\sigma(o^a) - y) \sigma(o^a)(1 - \sigma(o^a))}$$

f)

$$L_{CE}(\sigma(o^a), y) = -y \log(\sigma(o^a)) - (1-y) \log(1-\sigma(o^a))$$

where $\sigma(o^a) \in R$, $y \in \{0, 1\}$, $L_{CE} \in R$

thus $\frac{\partial L_{CE}}{\partial o^a} \in R$

$$\frac{\partial L_{CE}}{\partial o^a} = \frac{-y}{\sigma(o^a)} \cdot \frac{\partial \sigma(o^a)}{\partial o^a} - \frac{(1-y)}{1-\sigma(o^a)} \cdot \frac{\partial (1-\sigma(o^a))}{\partial o^a}$$

since $\frac{\partial \sigma(o^a)}{\partial o^a} = \sigma(o^a)(1-\sigma(o^a))$

Plugging it into above expression we get,

$$\frac{\partial L_{CE}}{\partial o^a} = -y(1-\sigma(o^a)) + (1-y)\sigma(o^a)$$

$$\boxed{\frac{\partial L_{CE}}{\partial o^a} = \sigma(o^a) - y}$$

$$\frac{\partial L_{CE}}{\partial o^a} \in R$$

$$g) L_{CE} = -y \log(\sigma(a)) - (1-y) \log(1-\sigma(a))$$

$$\frac{\partial L_{CE}}{\partial a} = \sigma(a) - y$$

$$L_{MSE} = (\sigma(a) - y)^2 \quad \frac{\partial L_{MSE}}{\partial a} = 2(\sigma(a) - y)\sigma(a)(1-\sigma(a))$$

Why cross entropy is better loss function than MSE?

- 1) Loss penalties are higher in CE for misclassification as compared to MSE. We can see this with an example:

$$\text{let } \sigma(a) = [0.1, 0.9]$$

$$y = [1, 0]$$

$$L_{CE} = -\log(0.1) = 1 \quad \checkmark$$

$$L_{MSE} = (0.1 - 1)^2 = 0.81$$

A larger value of loss ensures that the gradients computed are not squashed to zero particularly for deep networks. We can see this below.

$$\begin{aligned} \frac{\partial L_{CE}}{\partial a} &= \sigma(a) - y = -0.8 \\ \frac{\partial L_{MSE}}{\partial a} &= 2(\sigma(a) - y)\sigma(a)(1-\sigma(a)) = -0.162 \end{aligned} \quad \left. \begin{array}{l} \text{gradients} \\ \text{get squashed} \\ \text{in MSE loss} \\ \text{as compared to} \\ \text{CE} \end{array} \right\}$$

$$2) \text{Convex objective.} \rightarrow \frac{\partial^2 L_{CE}}{\partial a^2} = \sigma(a)(1-\sigma(a)) > 0 \quad \forall a \in \mathbb{R}$$

Convex objectives are easier to minimize thus its preferred loss function over MSE

$$h) L(x, y) = -\log o_y^s(x) ; x \in R^d , y \in R$$

For softmax activation, o_y^s is defined as:

$$o_y^s = \frac{e^{o_y^a}}{\sum_{k=1}^m e^{o_k^a}}$$

Thus we get,

$$\begin{aligned} L(x, y) &= -\log \left(\frac{e^{o_y^a}}{\sum_{k=1}^m e^{o_k^a}} \right) \\ &= -\log \left(e^{o_y^a} \right) + \log \left(\sum_{k=1}^m e^{o_k^a} \right) \end{aligned}$$

$$L(x, y) = -o_y^a + \log \sum_{k=1}^m e^{o_k^a}$$

$$i) \hat{R} = \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}^i, \mathbf{y}^i)$$

Set of parameters θ associated are:

$$w^{(1)} \text{ Dimension} = d_h \times d$$

$$b^{(1)} \text{ Dimension} = d_h \times 1$$

$$w^{(2)} \text{ Dimension} = m \times d_h$$

$$b^{(2)} \text{ Dimension} = m \times 1$$

$$\text{Total scalar parameter } n_\theta = d_h(d+1) + m(d_h+1)$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \hat{R}$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}^i, \mathbf{y}^i)$$

$$j) \quad \theta_{t+1} = \theta_t - \alpha \frac{\partial \hat{R}}{\partial \theta}$$

$$\boxed{\theta_{t+1} = \theta_t - \frac{\alpha}{n} \sum_{i=1}^n \frac{\partial L(x^i, y^i)}{\partial \theta}}$$

n = size of training dataset

α = learning rate ; $\alpha > 0$

$$k) \quad \hat{R}_{reg}(\theta) = \hat{R}(\theta) + \lambda \|\theta\|_2^2$$

$$\frac{\partial \hat{R}_{reg}(\theta)}{\partial \theta} = \frac{\partial \hat{R}(\theta)}{\partial \theta} + 2\lambda \|\theta\|_2$$

plugging this back in gradient descent equation

$$\boxed{\theta_{t+1} = \theta_t - \alpha \left[\frac{1}{n} \sum_{i=1}^n \frac{\partial L(x^i, y^i)}{\partial \theta} + 2\lambda \|\theta\|_2 \right]}$$

$$l) \quad L(x, y) = -\log o_y^s$$

From the expression of loss derived in part (h), we get

$$L(x, y) = -o_y^a + \log \sum_{k=1}^m e^{o_k^a}$$

$$\frac{\partial L}{\partial o_k^a} \in R^m$$

Computing partial derivative wrt any 'k' we get

$$\frac{\partial L}{\partial o_{k=y}^a} = -1 + \frac{1}{\sum_{j=1}^m e^{o_j^a}} (0 + \dots e^{o_y^a} + \dots 0)$$

$$\frac{\partial L}{\partial o_{k=y}^a} = \frac{e^{o_y^a}}{\sum_{j=1}^m e^{o_j^a}} - 1 = o_y^s - 1$$

$$\text{for } k \neq y, \quad \frac{\partial o_y^a}{\partial o_{k \neq y}^a} = 0$$

hence,

$$\frac{\partial L}{\partial o_k^a} = 0 + \frac{1}{\sum_{j=1}^m e^{o_j^a}} \left(e^{o_k^a} \right)$$

$$= o_k^s$$

therefore, $\frac{\partial L}{\partial o_k^a} = \begin{cases} o_k^s - 1 & ; k = y \\ o_k^s & ; k \neq y \end{cases}$

In vector form we can thus write

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

where $\text{onehot}_m(y) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \leftarrow y^{\text{th}} \text{ index}$

$$\frac{\partial L}{\partial o^a} \in \mathbb{R}^m ; o^s \in \mathbb{R}^m$$

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

From part-(l) we know that

$$\frac{\partial L}{\partial o_k^a} = \begin{cases} o_k^s - 1 & ; \quad k=y \\ o_k^s & ; \quad k \neq y \end{cases}$$

$$o^a = w^{(2)} h^s + b^{(2)} ; \quad w^{(2)} \in R^{m \times d_h}$$

$$b^{(2)} \in R^m$$

$$h^s \in R^{d_h}$$

$$o_k^a = \sum_{j=1}^{d_h} w_{kj}^{(2)} h_j^s + b_k$$

$$\text{thus, } \frac{\partial o_k^a}{\partial w_{kj}^{(2)}} = h_j^s$$

$$\frac{\partial o_k^a}{\partial b_k} = 1$$

substituting the values derived we get

$$\frac{\partial L}{\partial w_{kj}^{(2)}} = \begin{cases} (o_k^s - 1) h_j^s & ; k = y \\ o_k^s h_j^s & ; k \neq y \end{cases}$$

$$E \quad k \in \{1, 2, \dots, m\} \quad \& \quad j \in \{1, 2, \dots, d_n\}$$

$$\frac{\partial L}{\partial b_k^{(2)}} = \begin{cases} (o_k^s - 1) & k = y \\ o_k^s & k \neq y \end{cases}$$

$$n) \quad \frac{\partial L}{\partial w^{(2)}} = \frac{\partial L}{\partial o^a} \frac{\partial o^a}{\partial w^{(2)}}$$

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

$$L \in \mathbb{R} \quad o^a \in \mathbb{R}^m \Rightarrow \frac{\partial L}{\partial o^a} \in \mathbb{R}^m$$

$$w^{(2)} \in \mathbb{R}^{m \times d_h}, \quad o^a \in \mathbb{R}^m$$

$\frac{\partial o^a}{\partial w^{(2)}}$ is a tensor with dim = $m \times m \times d_h$

where $\frac{\partial o^a}{\partial w^{(2)}} = \begin{bmatrix} \frac{\partial o_1^a}{\partial w^{(2)}} \\ \vdots \\ \frac{\partial o_m^a}{\partial w^{(2)}} \end{bmatrix}_{m \times 1}$

for any $k \in \{1, 2, \dots, m\}$ $\frac{\partial o_k^a}{\partial w^{(2)}} \in \mathbb{R}^{m \times d_h} = \begin{bmatrix} o^T \\ \vdots \\ o^{s^T} \\ \vdots \\ o^T \end{bmatrix}_{m \times d_h}; \quad o^T \in \mathbb{R}^{d_h}$

$\leftarrow k^{\text{th}} \text{ index}$

$$\text{where, } h^s = \begin{bmatrix} h_1^s & h_2^s & \dots & h_{d_h}^s \end{bmatrix}^T \in \mathbb{R}^{d_h} \quad \text{and } o^T = \begin{bmatrix} 0 & 0 & \dots & 0 \end{bmatrix}^T \in \mathbb{R}^{d_h}$$

Combining we get

$$\frac{\partial L}{\partial w^{(2)}} \quad \text{dimension} = m \times d_h$$

$$\frac{\partial L}{\partial o} \in \mathbb{R}^m$$

$$\frac{\partial o}{\partial w^{(2)}} \quad \text{dimension} = m \times m \times d_h$$

The vector-tensor product upon multiplication is equivalent to the below matrix product

$$\frac{\partial L}{\partial w^{(2)}} = \begin{bmatrix} o_1^s \\ \vdots \\ o_g^s - 1 \\ \vdots \\ o_m^s \end{bmatrix}_{m \times 1} \begin{bmatrix} h_1^s & h_2^s & \dots & h_{d_h}^s \end{bmatrix}^T_{1 \times m}$$

$$\boxed{\frac{\partial L}{\partial w^{(2)}} = \left(o^s - \text{onehot}_m(y) \right)_{m \times 1} h^s_{1 \times m}^T}$$

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

$$\frac{\partial L}{\partial o^a} \text{ dimension} = 1 \times m, \quad b^{(2)} \in \mathbb{R}^m$$

$$o^a \in \mathbb{R}^m \Rightarrow \frac{\partial o^a}{\partial b^{(2)}} \text{ dimension} = m \times m$$

we know

$$\frac{\partial o_k^a}{\partial b_k^2} = 1$$

therefore

$$\frac{\partial o^a}{\partial b^{(2)}} = I_{m \times m} = \underbrace{\begin{bmatrix} 1 & 0 & \dots & \dots & 0 \\ 0 & 1 & \dots & & 0 \\ \vdots & \ddots & \dots & 2 & \vdots \\ 0 & 0 & \dots & & 1 \end{bmatrix}}$$

$$\frac{\partial L}{\partial b^2} = \left(o^s - \text{onehot}_m(y) \right)^T I_{m \times m}$$

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

$$\frac{\partial L}{\partial o_k^a} = \begin{cases} o_k^s - 1 & ; \quad k=y \\ o_k^s & ; \quad k \neq y \end{cases}$$

$$o_k^a = \sum_{j=1}^{d_h} w_{kj}^2 h_j^s + b_k^{(2)}$$

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

therefore,

$$\frac{\partial L}{\partial h_j^s} = \sum_{k=1}^m o_{k \neq y}^s w_{kj}^{(2)} + (o_y^s - 1) w_{yj}^{(2)}$$

b) Dimensions

$$\frac{\partial L}{\partial h^s} = 1 \times d_h$$

$$\frac{\partial L}{\partial o^a} = 1 \times m$$

$$\frac{\partial o^a}{\partial h^s} = m \times d_h$$

$$\frac{\partial L}{\partial o^s_a} = o^s - \text{onehot}_m(y) \quad ; \quad o^s \in \mathbb{R}^m \quad \text{onehot}_m(y) \in \mathbb{R}^m$$

$$\frac{\partial o_k^a}{\partial h_j^s} = w_{kj}^{(2)} \quad \forall k \in \{1, 2, \dots, m\}, j \in \{1, 2, \dots, d_h\}$$

$$\frac{\partial o^a}{\partial h^s} = w^{(2)} \quad ; \quad w^{(2)} \in \mathbb{R}^{m \times d_h}$$

$$\boxed{\frac{\partial L}{\partial h^s} = \left(o^s - \text{onehot}_m(y) \right)^T w^{(2)}_{m \times d_h}}$$

$$2) \frac{\partial L}{\partial h_j^a} = \frac{\partial L}{\partial h_j^s} \cdot \frac{\partial h_j^s}{\partial h_j^a}$$

$$h_j^s = \ln(1 + \exp(h_j^a))$$

$$\frac{\partial h_j^s}{\partial h_j^a} = \frac{1}{1 + \exp(h_j^a)} \cdot \exp(h_j^a)$$

In general if $i, j \in \{1, 2, \dots, d_n\}$

$$\frac{\partial h_i^s}{\partial h_j^a} = \begin{cases} \frac{1}{1 + \exp(-h_j^a)} & i = j \\ 0 & i \neq j \end{cases}$$

from part-(o) we derived earlier,

$$\frac{\partial L}{\partial h_j^s} = \sum_{k=1}^m o_{k \neq j}^s w_{kj}^{(2)} + (o_j^s - 1) w_{jj}^{(2)}$$

therefore we get

$$\frac{\partial L}{\partial h_j^a} = \left(\sum_{k=1}^m o_{k \neq j}^s w_{kj}^{(2)} + (o_j^s - 1) w_{jj}^{(2)} \right) \frac{1}{1 + \exp(-h_j^a)}$$

$$r) \quad \frac{\partial L}{\partial h^a} = \frac{\partial L}{\partial h^s} \frac{\partial h^s}{\partial h^a}$$

$$\frac{\partial L}{\partial h^s} = \left(0^S - \text{onehot}_m(y) \right)^T_{1 \times m} W^{(2)}_{m \times d_h}$$

from
part-(b)

Dimensions: $\frac{\partial L}{\partial h^s} = 1 \times d_h$

$$W^{(2)} \rightarrow m \times d_h$$

$$\frac{\partial h_i^s}{\partial h_j^a} = \begin{cases} \frac{1}{1 + \exp(-h_j^a)} & i = j \\ 0 & i \neq j \end{cases}$$

$$\frac{\partial h^s}{\partial h^a} = \text{diag}(H)$$

$$\text{where } H = \begin{bmatrix} \frac{1}{1 + e^{-h_1^a}} & \frac{1}{1 + e^{-h_2^a}} & \dots & \frac{1}{1 + e^{-h_{d_h}^a}} \end{bmatrix}_{d_h \times 1}^T$$

Dimensions

$$\frac{\partial h^s}{\partial h^a} = d_h \times d_h$$

$$h^a, h^s \in \mathbb{R}^{d_h}$$

$$\frac{\partial L}{\partial h^a} = \left(o^s - \text{onehot}_m(y) \right)^T_{1 \times d_h} W^{(2)}_{m \times d_h} \text{diag}(H)_{d_h \times d_h}$$

$$\frac{\partial L}{\partial h^a} \quad \text{dimension} = 1 \times d_h$$

IFT6390 HW3 Theory Q3

$$3. \text{ a) } h = \left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1,$$

where i = input height

p = padding,

k = kernel size,

s = stride

Let h_1 , h_2 and h_3 be the height of the first, second and third layers respectively.

$$h_1 = \left\lfloor \frac{128 + 2 \times 3 - 8}{2} \right\rfloor + 1$$

$$= 64$$

Since input image is 128×128 , and 32 kernels were convolved over it, the dimensions of the first layer is $32 \times 64 \times 64$.

Since the second layer is a 2×2 max-pooling with stride 2,

$$h_2 = \left\lfloor \frac{64}{2} \right\rfloor = 32.$$

\Rightarrow The dimensions of second layer is $32 \times 32 \times 32$.

Since third layer convolves 64 3×3 kernels with stride 1 and zero-padding of 1,

$$h_3 = \left\lfloor \frac{32 + 2 \times 1 - 3}{1} \right\rfloor + 1$$

$$= 32$$

\Rightarrow The dimensions of third layer is $64 \times 32 \times 32$.
(last)

$$3b) . \text{ Number of parameters} = 64 \times 3 \times 3 \times 32 \\ = 18432$$

3c) Let x be $[1, 1, 4, 4, 4, 1, 1]$
k be $[1/4, 1/4, 1/4]$

Let x^{full} which is x with padding for full convolution,
be $[0, 0, 1, 1, 4, 4, 4, 1, 1, 0, 0]$

Full convolution:

$$(x^{\text{full}} * k)_i = \sum_p x_{i+p}^{\text{full}} k_p$$

$$(x^{\text{full}} * k) = [0.25, 0.5, 1.5, 2.25, 3, 2.25, 1.5, 0.5, 0.25]$$

Valid convolution:

Let x^{valid} be x with padding for valid convolution,
which is simply x (no padding).

$$(x^{\text{valid}} * k)_i = \sum_p x_{i+p}^{\text{valid}} k_p$$

$$(x^{\text{valid}} * k) = [1.5, 2.25, 3, 2.25, 1.5]$$

Same convolution:

Let x^{same} be $[0, 1, 1, 4, 4, 4, 1, 1, 0]$

$$(x^{\text{same}} * k)_i = \sum_p x_{i+p}^{\text{same}} k_p$$

$$(x^{\text{same}} * k) = [0.5, 1.5, 2.25, 3, 2.25, 1.5, 0.5]$$

3d) The convolution operation effectively computes the average of the values in x within the window of the kernel. Since the kernel size is 3, every 3 consecutive values in x are averaged. It is an average pooling convolution.