

Q3.

a) 
$$\frac{(\text{input-size} + 2 \times \text{padding-size} - \text{kernel-size})}{\text{stride}}$$

First layers

$$\frac{(128 + 6 - 2)}{2} = \frac{128}{2} = 64$$

$$\Rightarrow \boxed{64 \times 64 \times 32}$$

Second layer: (max pooling & non-overlapping)

$$\frac{64}{2} = 32 \Rightarrow 32 \times 32 \times 32$$

after  
max  
Pooling

Third layers (64 3x3, stride=1 zero-padding=1)

$$\Rightarrow \frac{32 + 2 - 3}{1} + 1 = 32$$

$$\Rightarrow \boxed{64 \times 32 \times 32}$$

b)

Number of parameters needed in the Conv layer  
without considering bias

$$(m \times n \times d) \times K \Rightarrow 3 \times 3 \times 32 \times 64 = 18432$$

Filter height      num of filters       $\rightarrow$  num of  
width      in previous      filters in  
                layer      this layer

C)

Same Convolution

Putting Center of filter on the first element.

$$\begin{bmatrix} 0 & [1, 1, 4, 4, 4, 1, 1] & 0 \end{bmatrix}$$
$$[\frac{1}{4}, \frac{1}{4}, \frac{1}{4}]$$

↗ slide

$$= \left[ \frac{2}{4}, \frac{6}{4}, \frac{9}{4}, 3, \frac{9}{4}, \frac{6}{4}, \frac{3}{4} \right]$$

Full Convolution

Bottom Right element of filter on Top Left element  
of matrix.

$$\begin{bmatrix} 0, 0 & [1, 1, 4, 4, 4, 1, 1] & 0, 0 \end{bmatrix}$$
$$[\frac{1}{4}, \frac{1}{4}, \frac{1}{4}]$$

↗ slide

$$= \left[ \frac{1}{4}, \frac{2}{4}, \frac{6}{4}, \frac{9}{4}, 3, \frac{9}{4}, \frac{6}{4}, \frac{3}{4}, \frac{1}{4} \right]$$

Valid Convolution:

normal Convolution

$$[1, 1, 4, 4, 4, 1, 1]$$
$$[\frac{1}{4}, \frac{1}{4}, \frac{1}{4}]$$

↗ slide

$$= \left[ \frac{6}{4}, \frac{9}{4}, 3, \frac{9}{4}, \frac{6}{4} \right]$$

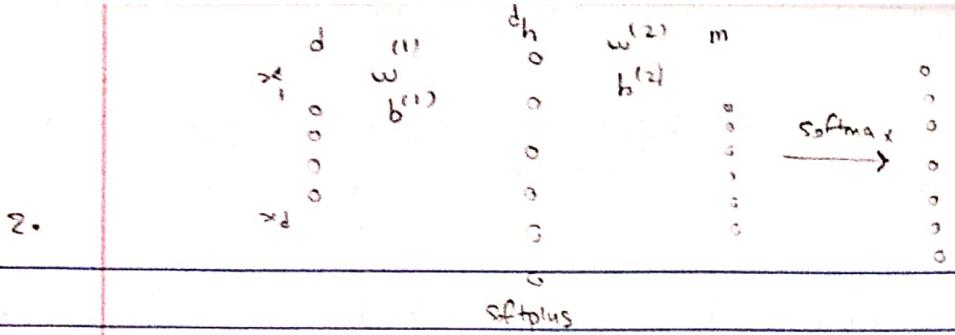
d)

- Valid & Full Convolution tend to keep data from the corner of the image while

Same Convolution uses no padding and might result in loss of data.

- Valid & Full Convolution makes the image bigger while same Convolution reduces the size.

- Same Convolution is practical because at the end we want lower dimension representation of our data.



a) -  $b^{(1)}$  contains a bias term for every  $d_h$  neurons in the hidden layer  
so it is a  $d_h \times 1$  vector

-  $h^a$  in matrix form is: 
$$h^a = w^{(1)} \cdot x + b^{(1)} \in \mathbb{R}^{d_h}$$

one element of  $h^a \rightarrow h_j^a$  jth row of  $w^{(1)} \rightarrow w_{j\cdot}$   
jth element of  $b^{(1)} \rightarrow b_j^{(1)}$

$$h_j^a = w_{j\cdot}^{(1)} \cdot x + b_j^{(1)} = \left( \sum_{i=1}^d w_{ji}^{(1)} \cdot x_i \right) + b_j^{(1)}$$

- output vector of the hidden layer  $h^s$  with respect to  $h^a$  is:

$$h^s = \text{softplus}(h^a) \rightarrow \text{softplus is applied on each } d_h \text{ elements of } h^a \rightarrow h^s \in \mathbb{R}^{d_h}$$

b)

$w^{(2)}$  → just like what we had in  $w^{(1)}$  each row of  $w^{(2)}$  represents the weights for every neuron in the output layer ( $m$ ) and each column represent weights for every neuron in the hidden layer ( $d_h$ )

So  $w^{(2)} \in \mathbb{R}^{m \times d_h}$

and likewise  $b^{(2)}$  is  $\mathbb{R}^m$  →  $m$  bias term per  $m$  output neurons

$$o^a = w^{(2)} \cdot h^s + b^{(2)} \in \mathbb{R}^m$$

one element of  $o^a \rightarrow o_k^a \Rightarrow$  kth element of  $w^{(2)} \rightarrow w_{k\cdot}$   
kth element of  $b^{(2)} \rightarrow b_k^{(2)}$

$$o_k^a = w_{k\cdot}^{(2)} \cdot h^s + b_k^{(2)} = \left( \sum_{i=1}^{d_h} w_{ki}^{(2)} \cdot h_i^s \right) + b_k^{(2)}$$

c)

$$o^s = \text{softmax}(o^a) = \frac{\exp(o^a)}{\sum_{k=1}^m \exp(o_k^a)}$$

$$o_k^s = \text{softmax}(o_k^a) = \frac{\exp(o_k^a)}{\sum_{i=1}^m \exp(o_i^a)}$$

Since  $\exp(x)$  is always positive,  $o_k^s$  is just division of some  $\exp(x)$  function so  $o_k^s$  are always positive.

$$\sum_{k=1}^m o_k^s = \sum_{k=1}^m \frac{\exp(o_k^a)}{\sum_{i=1}^m \exp(o_i^a)} = \frac{\sum_{k=1}^m \exp(o_k^a)}{\sum_{i=1}^m \exp(o_i^a)} = 1$$

why it is important? because the task is a classification problem and using a softmax as last layer has the interpretation of probability, if the test input belongs to each of the  $m$  class so according to the laws of probability it should sum to one to be a valid probability distribution.

d)

$$\text{softmax} \Rightarrow p(Y=k | X) = \frac{e^{o_k^a}}{\sum_{i=1}^m e^{o_i^a}} = \frac{e^{o_k^a}}{\sum_{i=1}^m (o_i^a - o_k^a)}$$

$$= \frac{1}{\sum_{i=1}^m e^{-(o_k^a - o_i^a)}} \quad \text{if } m=2 \rightarrow \text{binary classification}$$

assume labels are  $\{0, 1\}$

(2)

$$P(Y=1|X) = \frac{1}{e^{-(\theta_1^a - \theta_2^a)} + e^{-(\theta_1^a - \theta_2^a)}} = \frac{1}{e^{-(\theta_1^a - \theta_2^a)} + 1}$$

tion

$$= \frac{1}{1 + e^{-\theta^a}} \quad (\text{where } \theta^a = \theta_1^a - \theta_2^a)$$

$$= \delta(\theta^a) \text{ - sigmoid}$$

e)

$$L_{MSE}(\delta(\theta^a), Y) \quad \delta(\theta^a) = \sigma^{(s)}$$

$$\frac{\partial L_{MSE}(\theta^a, Y)}{\partial \theta^a} = \frac{\partial L(\theta^a, Y)}{\partial \theta^a} \times \frac{\partial \theta^a}{\partial \theta^a}$$

$$= \boxed{(\delta(\theta^a) - Y) \times \delta(\theta^a)(1 - \delta(\theta^a))}$$

$$f) L_{CE}(\delta(\theta^a), y) = -(y \log(\delta(\theta^a)) + (1-y) \log(1 - \delta(\theta^a)))$$

$n$  = # of training example

$y$  = true value

$\delta(\theta^a)$  = predicted value

$$\frac{\partial L_{CE}(\delta(\theta^a), y)}{\partial \theta^a} = \frac{\partial L_{CE}(\delta(\theta^a), y)}{\partial \delta(\theta^a)} \times \frac{\partial \delta(\theta^a)}{\partial \theta^a}$$

$$\frac{\partial L_{CE}(\delta(\theta^a), y)}{\partial \delta(\theta^a)} = \boxed{-\left(\frac{y}{\delta(\theta^a)} - \frac{(1-y)}{1-\delta(\theta^a)}\right) \times \delta(\theta^a)(1 - \delta(\theta^a))}$$

g) first using MSE means that we assume that the underlying data has been generated from a normal distribution in bayesian term this means that we assume a gaussian prior while in reality a dataset that can be classified into 2 categories is not from a normal distribution but a bernoulli distribution

Secondly → MSE Function is non convex for binary classification (why?) in other way it is not guaranteed to minimize cost function - this because MSE expects real valued input in range  $(-\infty, +\infty)$  while binary classification model output probabilities in range  $(0,1)$  through the sigmoid function.

(3)

3) so when a bounded value from a sigmoid function is passed to the MSE  
the result is not convex. on one side the function is concave while on the other side the function  
is convex and no clear minimum point

in addition binary cross entropy loss is equivalent to fitting the model using  
maximum likelihood estimation

and the decision boundary in a classification task is large (in comparison with regression)  
but MSE doesn't punish missclassification enough.

h)

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

$$= -\log \text{softmax}(o_y^s(x))$$

$$= -\log \frac{\exp(o_y^s(x))}{\sum_{i=1}^m \exp(o_i^s(x))}$$

i)

$$\hat{R}(f, D) = \frac{1}{n} \sum_{i=1}^n L(f_\theta(x^{(i)}), y^{(i)})$$

which is the average over a finite dataset D

$\theta$  is the parameter that the model learns during training, include all the weight and biases in all the layers:

$$\theta = [w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}]$$

$w^{(1)} \in \mathbb{R}^{d_h \times d}$     $b^{(1)} \in \mathbb{R}^{d_h}$     $w^{(2)} \in \mathbb{R}^{m \times d_h}$     $b^{(2)} \in \mathbb{R}^m$

$$n_\theta = d_h \times (d+1) + m \times (d_h+1) \quad \text{scalar parameters}$$

The optimization problem of training the model to find the optimal values of parameter is:

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

j) learning rate

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} \hat{R} = \theta^{(t)} - \eta \frac{\partial}{\partial \theta} \left( \frac{1}{n} \sum_{i=1}^n L(f_{\theta}(x^{(i)}), y^{(i)}) \right)$$

$$S = [w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}]$$

$$w^{(1)} = w^{(1)} - \eta \frac{\partial}{\partial w^{(1)}} \hat{R}$$

$$w^{(2)} = w^{(2)} - \eta \frac{\partial}{\partial w^{(2)}} \hat{R}$$

$$b^{(1)} = b^{(1)} - \eta \frac{\partial}{\partial b^{(1)}} \hat{R}$$

$$b^{(2)} = b^{(2)} - \eta \frac{\partial}{\partial b^{(2)}} \hat{R}$$

k)

$$\theta^{(t+1)} = \theta^{(t)} - \eta \frac{\partial}{\partial \theta} \left( \frac{1}{n} \sum_{i=1}^n L(f_{\theta}(x^{(i)}), y^{(i)}) \right) - \eta \frac{\partial}{\partial \theta} (\lambda \|\theta\|_2^2)$$

l)

$$L(x, y) = -\log o_y^S(x) = -\log \text{softmax}(o_y^a(x)) = -\log \frac{\exp(o_y^a(x))}{\sum_{i=1}^m \exp(o_i^a(x))}$$

$$\frac{\partial L(x, y)}{\partial o_k^a} = \frac{\sum_{i=1}^m \exp(o_i^a(x))}{\exp(o_y^a(x))} x \frac{\partial}{\partial o_k^a} \left( \frac{\exp(o_y^a(x))}{\sum_{i=1}^m \exp(o_i^a(x))} \right)$$

for  $k \neq y$

$$+ \frac{\exp(o_k^a(x)) \exp(o_y^a(x))}{(\sum_{i=1}^m \exp(o_i^a(x)))^2} x \frac{\exp(o_y^a(x))}{\exp(o_k^a(x))}$$

$$= \frac{\exp(o_k^a(x))}{\sum_{i=1}^m \exp(o_i^a(x))} = o_k^S$$

(5)

$$\frac{\partial L(x, y)}{\partial o_y^a} = \frac{-\sum_{i=1}^m \exp(o_i^a(x))}{\exp(o_y^a(x))} \times \frac{\partial}{\partial o_y^a} \left( \frac{\exp(o_y^a(x))}{\sum_{i=1}^m \exp(o_i^a(x))} \right)$$

$$= - \frac{\exp(o_y^a(x)) \sum_{i=1}^m \exp(o_i^a(x)) - \exp(o_y^a(x))^2 \sum_{i=1}^m \exp(o_i^a(x))}{\left( \sum_{i=1}^m \exp(o_i^a(x)) \right)^2} \times \frac{\exp(o_y^a(x))}{\exp(o_y^a(x))}$$

$$= - \frac{\exp(o_y^a(x)) \left( \sum_{i=1}^m \exp(o_i^a(x)) - \exp(o_y^a(x)) \right)}{\sum_{i=1}^m \exp(o_i^a(x))} \times \frac{1}{\exp(o_y^a(x))}$$

$$\Rightarrow - \left[ 1 - \frac{\exp(o_y^a(x))}{\sum_{i=1}^m \exp(o_i^a(x))} \right] = -1 + o_y^s(x)$$

$$\frac{\partial L(x, y)}{\partial o_y^a} = \begin{bmatrix} o_1^s & \dots & o_m^s & \dots & 0 \\ \vdots & & \vdots & & \vdots \\ o_y^s & -1 & o_y^s(x) & \dots & 0 \\ \vdots & & \vdots & & \vdots \\ o_m^s & \dots & o_m^s(x) & \dots & 0 \end{bmatrix} = 0 \quad \text{--- onehot } \begin{matrix} s \\ m \end{matrix}$$

(S)

(6a)

(6)

m)

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

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

$$o_k^a = w_{kj,:}^{(2)} \cdot h^{(S)} + b_k^{(2)} \Rightarrow \frac{\partial o_k^a}{\partial w_{kj}^{(2)}} = \frac{\partial w_{kj,:}^{(2)} \cdot h^{(S)}}{\partial w_{kj}^{(2)}} + \frac{\partial b_k^{(2)}}{\partial w_{kj}^{(2)}} = \boxed{h_j^{(S)}}$$

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

$$\frac{\partial L(x,y)}{\partial o_k^a} = o_k^S(x) - \text{onehot}_m(y) \quad \text{from part (L)}$$

$$\rightarrow \boxed{\frac{\partial L(x,y)}{\partial w_{kj}^{(2)}} = \frac{\partial L(x,y)}{\partial o_k^a} \times h_j^{(S)} = (o_k^S(x) - \underbrace{\text{onehot}_m(y)}_{k\text{th entry of } L} ) \times h_j^{(S)}}$$

↳ one-hot representation of label.

$$\boxed{\frac{\partial L(x,y)}{\partial b_k^{(2)}} = \frac{\partial L(x,y)}{\partial o_k^a} \cdot 1 = \frac{\partial L(x,y)}{\partial o_k^a} = (o_k^S(x) - \text{onehot}_m(y))}$$

n)

$$\boxed{\frac{\partial L}{\partial w^{(2)}} = \frac{\partial L}{\partial o^a} \times \frac{\partial o^a}{\partial w^{(2)}} = (o^{(S)} - \text{onehot}_m(y)) \cdot h^{S^T}}$$

$$\boxed{\frac{\partial L}{\partial b^{(2)}} = \frac{\partial L}{\partial o^a} \times \frac{\partial o^a}{\partial b^{(2)}} = \frac{\partial L}{\partial o^a} = (o^{(S)} - \text{onehot}_m(y))}$$

$$\left( \frac{\partial L}{\partial o^a} = \underbrace{(o^{(S)} - \text{onehot}_m(y))}_{m \times 1} \right) \underbrace{h^{S^T}}_{1 \times d_h} \Rightarrow \boxed{\frac{\partial L}{\partial w^{(2)}} \text{ has dimension } m \times d_h \text{ like } w^{(2)}}$$

$$\boxed{\frac{\partial L}{\partial b^{(2)}} = \frac{\partial L}{\partial o^a} \Rightarrow \text{has dimension } m \times 1 \text{ like } b^{(2)}}$$

(7)

Q)

$$\frac{\partial L}{\partial h_j^S} = \sum_{k=1}^m \frac{\partial L}{\partial o_k^S} \frac{\partial o_k^S}{\partial h_j^S}$$

$$o_k^S = w_{k,j}^{(2)} \cdot h_j^{(S)} + b_{k,j}^{(2)} = \left( \sum_{j=1}^{d_h} w_{k,j}^{(2)} \cdot h_j^{(S)} \right) + b_{k,j}^{(2)}$$

$$\boxed{\frac{\partial o_k^S}{\partial h_j^{(S)}} = w_{k,j}^{(2)}}$$

$$\boxed{\frac{\partial L}{\partial o_k^S} = (o_k^S(x) - \text{onehot}_m(y))}$$

$$\boxed{\frac{\partial L}{\partial h_j^S} = \sum_{k=1}^m (o_k^S(x) - \text{onehot}_m(y))_k w_{k,j}^{(2)}}$$

P)

$$\frac{\partial L(x,y)}{\partial h_j^S} = \begin{bmatrix} \sum_{k=1}^m (o_k^S(x) - I(k=y)) w_{k,j}^{(2)} \\ \vdots \\ \sum_{k=1}^m (o_k^S(x) - I(k=y)) w_{k,d_h}^{(2)} \end{bmatrix}_{d_h \times 1}$$

$$= \begin{bmatrix} \sum_{k=1}^m o_k^S(x) w_{k,1}^{(2)} \\ \vdots \\ \sum_{k=1}^m o_k^S(x) w_{k,d_h}^{(2)} \end{bmatrix} - \begin{bmatrix} \sum_{k=1}^m I(k=y) w_{k,1}^{(2)} \\ \vdots \\ \sum_{k=1}^m I(k=y) w_{k,d_h}^{(2)} \end{bmatrix}$$

$$= w^{(2)T} o - w^{(2)T} \text{onehot}_m(y) = \boxed{w^{(2)T} (o^S - \text{onehot}_m(y))}$$

(8)

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

$$\text{onehot}_m(y) \rightarrow m \times 1 \Rightarrow \frac{\partial L(x, y)}{\partial h^{(1)}} = w^{(2)^T} (o^{(1)} - \text{onehot}_m(y))$$

$$o^{(1)} \rightarrow m \times 1$$

dimension of previous vector

$$\left( \frac{\partial L}{\partial h^{(1)}} \right)$$

$d_h \times m$

$m \times 1$

$$d_h \times 1$$

$$q) \quad \frac{\partial L}{\partial h_j^a} = \frac{\partial L}{\partial h_j^{(1)}} \times \frac{\partial h_j^{(1)}}{\partial h_j^a} \quad h_j^{(1)} = \text{softplus}(h_j^a)$$

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

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

$$\boxed{\frac{\partial L}{\partial h_j^a} = \left( \sum_{k=1}^m (o_k^{(1)}(x) - \text{onehot}_m(y)_k) w_{kj}^{(2)} \right) \frac{1}{1 + \exp(-h_j^a)}}$$

$$r) \quad \left[ \left( \sum_{k=1}^m (o_k^{(1)}(x) - I(k=y)) w_{kj}^{(2)} \right) \frac{1}{1 + \exp(-h_j^a)} \right]$$

$$\frac{\partial L}{\partial h_j^a} = \left[ \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right]$$

$$\left[ \left( \sum_{k=1}^m (o_k^{(1)}(x) - I(k=y)) w_{kj}^{(2)} \right) \frac{1}{1 + \exp(-h_j^a)} \right]$$

(g)

$$\boxed{w^{(2)T} \left( o^S - \text{onehot}_m(y_1) \right) \cdot \frac{1}{1 + \exp(-h^a)} = \frac{\partial L}{\partial h^a}}$$

$$\begin{pmatrix} o^S & \text{onehot}_m(y_1) \\ \downarrow & \downarrow \\ m \times 1 & m \times 1 \end{pmatrix} \rightarrow m \times 1$$

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

$$\frac{\partial L}{\partial h^a} = (d_h \times m) \times (m \times 1) \times (d_h \times 1)$$

$$= \boxed{\frac{\partial L}{\partial h^a} \rightarrow \text{dimension } d_h}$$

### HW 3

1-a)  $\text{Sign}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$

b)  $g(x) = \begin{cases} x & x > 0 \\ 0 & x < 0 \\ \text{not defined} & x=0 \end{cases} \Rightarrow g(x) = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \end{cases}$

at  $x=0$ ,  $g'(x)$  doesn't exist because of the discontinuity.

$g'(x^-)$  and  $g'(x^+)$  are not equal. So  $g'(x) = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \end{cases}$

$H(x) = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \\ \frac{1}{2} & x=0 \end{cases}$  so  $g'(x)$  is equal to  $H(x)$  where

$g'(x)$  exists (everywhere except  $x=0$ )

c) ①  $g(x) = xH(x)$

②  $g(x) = H(|x|) - H(-|x|)$

d)  $a(x) = \frac{1}{1 + e^{-kx}}$  when  $k$  is large  $\rightarrow$  we have 3 possible values for  $x$ :

①  $x > 0 \Rightarrow m = kx$  will be  $\infty$   $\text{as } x \rightarrow \infty \Rightarrow \lim_{m \rightarrow \infty} \frac{1}{1 + e^{-m}} = 1$  ✓  
because  $e^{-m}$  when  $m \rightarrow \infty$  will be  $e^{-\infty} \approx 0$

②  $x < 0 \Rightarrow m = kx$  will be  $-\infty$   $\text{as } x \rightarrow -\infty \Rightarrow \lim_{m \rightarrow -\infty} \frac{1}{1 + e^{-m}} = 0$  ✓  
because  $e^{-m}$  when  $m \rightarrow -\infty$  will be  $e^{(-\infty)} = e^{+\infty} = \infty$  and  $\frac{1}{\infty} = 0$

$$\textcircled{3} \quad x=0 \Rightarrow m=kx = k \cdot 0 = 0 \Rightarrow \frac{1}{1+e^0} = \frac{1}{2} \checkmark$$

so it is like  $H(n)$ .

$$e) s(x) = \frac{1}{1+e^{-x}} \Rightarrow \frac{ds(x)}{dx} = \frac{\frac{d}{dx}(1+e^{-x}) - \frac{d(1+e^{-x})}{dx}}{(1+e^{-x})^2}$$

$$\frac{0 - (1-1)(e^{-x})}{(1+e^{-x})^2} = \frac{e^{-x}}{1+e^{-x}} \times \frac{1}{1+e^{-x}} = \frac{e^{-x}}{1+e^{-x}} \times s(x) =$$

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

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \text{ where } \frac{\partial f_i}{\partial x_j} = \delta_{ij} s'(x_j) = \delta_{ij} (s(x_j) - s^2(x_j))$$

$$f) s(x) = \frac{1}{1+e^{-x}} \Rightarrow \ln(s(x)) = \ln((1+e^{-x})^{-1}) = -\ln(1+e^{-x})$$

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

so  $\ln(s(x)) = -\text{softplus}(-x) \checkmark$

4.

$$\begin{aligned} g) \text{ softplus}(n) - \text{softplus}(-n) &= \ln(1+e^n) - \ln(1+e^{-n}) \\ &= \ln(\alpha(n)) + \ln(\alpha(-n)) \rightarrow \text{based on previous part} \\ &= \ln\left(\frac{\alpha(n)}{\alpha(-n)}\right) = \ln\left(\frac{\frac{1}{1+e^{-n}}}{\frac{1}{1+e^n}}\right) = \ln\left(\frac{1+e^n}{1+e^{-n}}\right) \\ &= \ln\left(\frac{e^n(1+e^{-n})}{1+e^{-n}}\right) = \ln(e^n) = n \checkmark \end{aligned}$$

$$\begin{aligned} h) s(n+c)_i &= \frac{e^{(n_i+c)}}{\sum_j e^{(n_j+c)}} = \frac{e^c e^{n_i}}{e^c \sum_j e^{n_j}} = \frac{e^{n_i}}{\sum_j e^{n_j}} \\ &= s(n)_i \checkmark \end{aligned}$$

$$i) s(cx)_i = \frac{e^{(cn)_i}}{\sum_j e^{(cn)_j}} \quad cn=t \rightarrow s(t)_i =$$

$$\frac{e^{t_i}}{\sum_j e^{t_j}} \quad \left( \begin{array}{l} \text{if } s(t)_i = s(n)_i \Rightarrow \frac{e^{n_i}}{\sum_j e^{n_j}} = \frac{e^{t_i}}{\sum_j e^{t_j}} \\ s(n)_i = \frac{e^{n_i}}{\sum_j e^{n_j}} \end{array} \right)$$

In this case, we should have  $t=n$  in order to have the above equation correct  $\Rightarrow t=cn \Rightarrow cn \stackrel{?}{=} n$  it is only true if  $c=1$ . So in general we can say that softmax is not invariant under multiplication

$$\text{i. cont.) } S(cx) = \frac{e^{cx_i}}{\sum_j e^{cx_j}} = \frac{(e^{nx_i})^c}{\sum_j (e^{nx_j})^c} = \left( \frac{\sum_j (e^{nx_j})^c}{\sum_j e^{nx_j}} \right)^{-1}$$

considering this part, when  $c > 1$  then this part

will be greater if  $n_i < n_j$  and will be smaller if  $n_j > n_i$

Finally the result of the sum will be dependent to how many  $n_j$  are smaller than  $n_i$  and how much is the difference, because for example  $\frac{1}{2} + 1.25 < (1.25)^2 + (1.25)^2$  but  $\frac{1}{2} + 1.1 > (\frac{1}{2})^2 + (1.1)^2$

The same things happens when  $c < 1$  but in the opposite way. Finally the result will be to the power of  $-1$  which makes the result opposite if it was larger. The multiplications make the results more spiky and we get smaller and vice versa.

$$\text{j) } \frac{\partial S(x)}{\partial x_j} = \left[ \frac{\partial S(x)_1}{\partial x_j}, \frac{\partial S(x)_2}{\partial x_j}, \dots, \frac{\partial S(x)_n}{\partial x_j} \right] \text{ will be more confident}$$

$$\rightarrow \frac{\partial S(x)_i}{\partial x_j} = \frac{\partial}{\partial x_j} \left( \frac{e^{nx_i}}{\sum_t e^{nx_t}} \right) = \frac{-e^{nx_i} e^{nx_j}}{(\sum_t e^{nx_t})^2} = \frac{-x_i}{\sum_t e^{nx_t}} \frac{x_j}{\sum_t e^{nx_t}}$$

$$\text{if } i=j = -S(x)_i S(x)_j$$

$$\rightarrow \frac{\partial S(x)_i}{\partial x_j} = \frac{e^{nx_i} \sum_t e^{nx_t} - e^{nx_i} e^{nx_j}}{(\sum_t e^{nx_t})^2} = S(x)_i - S(x)_i S(x)_j$$

$\Rightarrow$  writing the above results together we get:

$$\frac{\partial S(x)_i}{\partial x_j} = \sum_{t=0}^n S(x)_t - S(x)_i S(x)_j = S(x)_i \left( \sum_{t=0}^n S(x)_t - S(x)_i \right)$$



b

K)  $\frac{\partial s(n)_i}{\partial n_j} = s(n)_i (\delta_{i,j} - s(n)_{ij})$  from previous part

$\Rightarrow$  the diagonal elements have  $i=j$  and will be:

$$s(n)_i = s(n)_i s(n)_j$$

and the others will be  $i \neq j \Rightarrow s(n)_i s(n)_j$

$$\Rightarrow \frac{\partial s(n)_i}{\partial n_j} = \text{diag}(s(n)) - \begin{bmatrix} s(n)_1 \\ s(n)_2 \\ \vdots \\ s(n)_n \end{bmatrix} [s(n)_1, s(n)_2, \dots, s(n)_n]$$

$$= \text{diag}(s(n)) - s(n) s(n)^T$$

m)

$$\nabla_u L(n, c) = \sum_{i=1}^k -c_i \nabla_n \log(s(n(u))_i) \quad \textcircled{1}$$

$$\frac{\partial \log(s(n(u))_k)}{\partial u_k} \stackrel{k=j}{=} \frac{s(n(u))_k - s(n(u))_k s(n(u))_j}{\log(s(n(u))_k)} \\ \xrightarrow{k \neq j} \frac{-s(n(u))_k s(n(u))_j}{\log(s(n(u))_k)}$$

using the result of previous question and equation \textcircled{1}

we get:

$$\nabla_u L(n, c) = -\sum_{i=1}^k c_i (\nabla_n n(u) - \nabla_u n(u); s(n(u))_i)$$

1)

L)

$$\nabla_u \log S(x(u))$$

$$= \frac{\partial \log S(x(u))}{\partial S(x(u))} \times \frac{\partial S(x(u))}{\partial x(u)} \times \frac{\partial x(u)}{\partial u}$$

$$= \frac{1}{S(x(u))} \times \left( \text{diag}(S(x(u)) - S(x(u))S(x(u))^T) \right) \nabla_x x(u)$$

$$= (I - S(x(u))^T) \nabla_u x(u)$$

$$= \nabla_u x(u) - A_u x(u) S(x(u))^T$$

$$\nabla_u x(u) - \underbrace{\sum_j \nabla_u x(u)_j S(x(u))_j^T}_{\sim}$$

according to the question  $S(x(u))_j$  is the probability of  $j$

$$= \nabla_u x(u) - E_j [\nabla_u x(u)_j]$$