## Q1.1

let $x'$ be $x + c$, $\quad x_i' = x_i + c$. ($x_i$ is value in $x'$, so does $x_i$).

$$\text{softmax}(x_i') = \frac{e^{x_i'}}{\sum_j e^{x_j'}} = \frac{e^{x_i + c}}{\sum_j e^{x_j + c}} = \frac{e^{x_i}}{\sum_j e^{x_j}} = \text{softmax}(x_i)$$

If we use $c = -\max x_i$, we can make sure $x_i' \leq 0$, and $e^{x_i'} \in [0, 1]$

In this way we can prevent overflow.

## Q1.2.

(1). the range of each element is $(0, 1)$, the sum is $1$.

(2). probability.

(3). The first step calcute the weight of each element, the second step sum the weights and the third step normalize the weights.

## Q1.3.

First we express a 2-layer NN into LR.

$$y = W_2(W_1 x + b_1) + b_2$$
$$\Leftrightarrow y = W_2 W_1 x + W_2 b_1 + b_2$$
$$\Leftrightarrow y = W' x + b' \quad (W' = W_2 W_1, \quad b' = W_2 b_1 + b_2).$$

We can recursively reduce a multi-layer NN to LR using the above way.

## Q1.4.

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

Q 1.5.

$$y = W^T x + b. \qquad y_j = \sum_{i=1}^{d} x_i W_{ij} + b_j$$
$$\underset{k \times 1}{} \ \underset{k \times d}{} \ \underset{d \times 1}{} \ \underset{k \times 1}{}.$$

I'll use the above notation.

① $\dfrac{dJ}{dW_{ij}} = \sum_{n=1}^{k} \left(\dfrac{dJ}{dy_n}\right)_* \times \left(\dfrac{dy_n}{dW_{ij}}\right) = \sum_{n=1}^{k} \delta_n \times \dfrac{d(\sum_{i=1}^{d} x_i W_{in} + b_n)}{dW_{ij}}$

$$= \delta_j \times \dfrac{d(\sum_{i=1}^{d} x_i W_{ij} + b_j)}{dW_{ij}} = \delta_j x_i$$

So $\dfrac{dJ}{dW} = \delta \cdot x^T$
$$\underset{k \times d}{} \qquad \underset{k \times 1}{} \ \underset{1 \times d}{}.$$

② $\dfrac{dJ}{dx_i} = \sum_{j=1}^{k} \left(\dfrac{dJ}{dy_j}\right) \times \left(\dfrac{dy_j}{dx_i}\right) = \sum_{j=1}^{k} \delta_j \times \dfrac{d\sum_{i=1}^{d} x_i W_{ij} + b_j}{dx_i}$

$$= \sum_{j=1}^{k} \delta_j W_{ij}$$

So $dJ/dx = (W\delta)^T$
$$\underset{1 \times d}{} \qquad \underset{d \times k, \ k \times 1}{}.$$

③ $\dfrac{dJ}{db_n} = \sum_{j=1}^{k} \left(\dfrac{dJ}{dy_j}\right) \times \left(\dfrac{dy_j}{db_n}\right) = \sum_{j=1}^{k} \delta_j \times \dfrac{d\sum_{i=1}^{d} x_i W_{ij} + b_j}{db_n} = \delta_n.$

So $dJ/db = \delta^T$
$$\underset{1 \times k}{} \qquad \underset{k \times 1}{}$$

Q1.6.

① take a 2-layer NN as example.

$z_1 = W_1 x + b_1$,  $y_1 = \sigma(z_1)$  $z_2 = W_2 y_1 + b_2$  $y_2 = \sigma(z_2)$

then  $\partial y_2 / \partial W_1 = \partial y_2 / \partial z_2 \cdot \partial z_2 / \partial y_1 \cdot \partial y_1 / \partial z_1 \cdot \partial z_1 / \partial W_1$

In the above equation $\partial y_2 / \partial z_2$ and $\partial y_1 / \partial z_1$ follows $\partial y / \partial z = \sigma(z)(1 - \sigma(z))$

When $z$ goes smaller, $\sigma(z)$ will be close to 0, making gradient close to 0.

② tanh $\in (-1, 1)$  sigmoid $\in (0, 1)$.

tanh has a better gradient when initializing around 0.

③ When initialized close to 0, tanh has a larger gradient than sigmoid.

④ $\tanh(x) = \dfrac{1 - e^{-2x}}{1 + e^{-2x}} = \dfrac{2}{1 + e^{-2x}} - 1 = 2\,\text{sigmoid}(2x) - 1.$

~~Note:~~