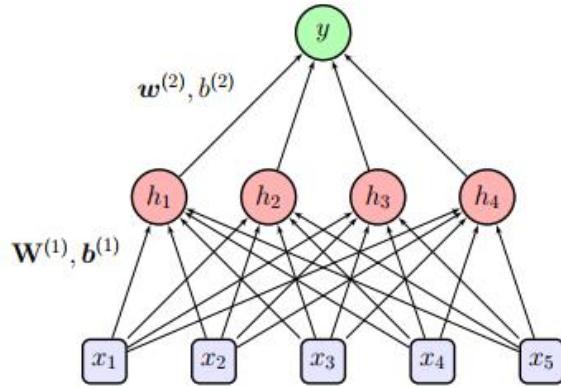


Homework Solns

Problem 1 [20 pts]: In this problem you will find a set of weights and biases for a multilayer perceptron which determine if a list of four numbers in the descending order. More specifically, you receive five inputs x_1, x_2, x_3, x_4 and x_5 , where $x_i \in R$, and the network must output 1 if $x_1 > x_2 > x_3 > x_4 > x_5$ and 0 otherwise. You will use the following Multilayer Perceptron (MLP) architecture consisting of one input layer with five nodes, one hidden layer with four nodes one output layer with one node.



The activation function of the hidden units and the output unit can be assumed to be a hard threshold function.

$$\sigma(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

Find a set of weights and biases for the network which correctly implements this function (including cases where some of the inputs are equal). Note, in some cases, if some of the inputs are equal to each other (say $x_2 = x_3$), your output should be 0. Your answer should include:

1. A 4×5 weight matrix $\mathbf{W}^{(1)}$ for the hidden layer.
2. A 4-dimensional vector of biases $\mathbf{b}^{(1)}$ for the hidden layer.
3. A 4-dimensional vector of weights $\mathbf{w}^{(2)}$ for the output layer.
4. A scalar bias $b^{(2)}$ for the output layer.

Solu: We need $x_1 > x_2 > x_3 > x_4 > x_5$, $x_i \in \text{real numbers}$

One way of ensuring $x_i > x_{i+1}$

is to ensure that $d_i = x_i - x_{i+1} > 0$, $i=1, 2, 3, 4$

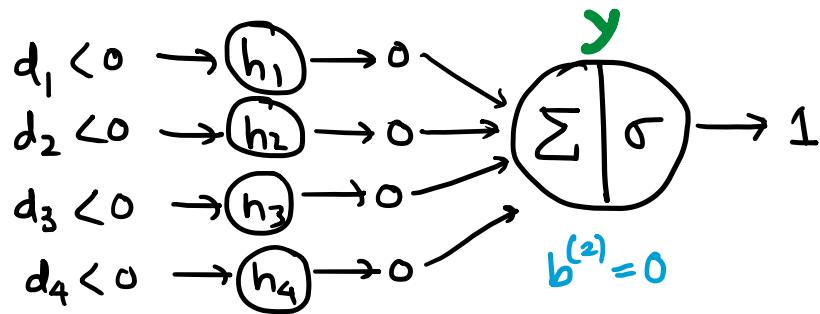
So we are given 4 perceptrons in the hidden layer, and each of them could be fed with one difference d_i . We could expect the perceptron h_i to fire only if $d_i > 0$, however, given our output activation fires even when $d_i = 0$, this way will not work when $x_i = x_{i+1}$ ($d_i = 0$).

So consider $d_i = x_{i+1} - x_i$. For this, when $d_i < 0$, the output would be zero, and for $d_i \geq 0$, the output would be 1

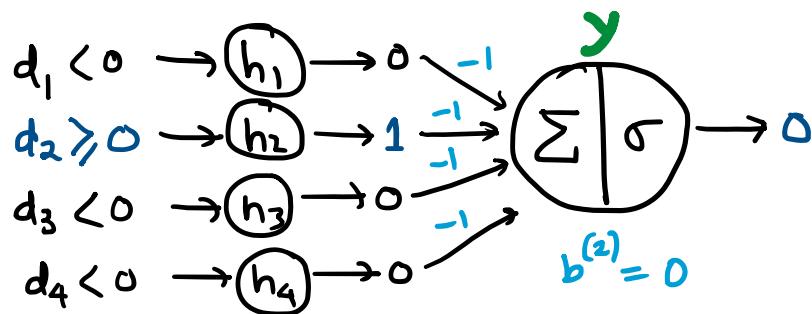
$$d_i > 0 \rightarrow h_i \rightarrow 1$$

$$d_i < 0 \rightarrow h_i \rightarrow 0$$

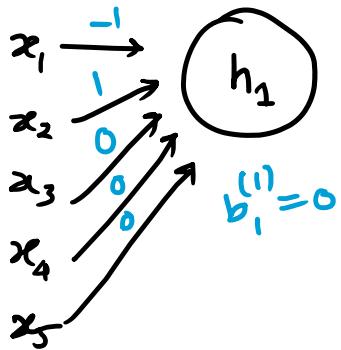
If all $d_i < 0$, $i=1,2,3,4$, then the output of hidden activations $h_i = 0 \forall i$, and hence the output would be set to 1



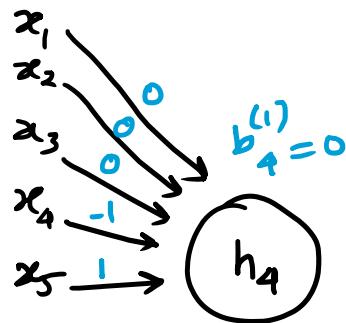
However, if any $d_i \geq 0$, then the output of the corresponding hidden unit would be 1, i.e. $h_i = 1$. To ensure that the output activation y does not fire when any of the h_i 's becomes 1, one could set the weights of the second layer to some negative value say $w_i^{(2)} = -1$



As for the weights of the first layer, we can define individual differences d_i as follows



...



(1) The weight matrix for the first layer is

$$\tilde{W}^{(1)} = \begin{bmatrix} -1 & 1 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 0 \\ 0 & 0 & -1 & 1 & 0 \\ 0 & 0 & 0 & -1 & 1 \end{bmatrix}$$

(2) The bias vector for first layer can be all zeros

$$\tilde{b}^{(1)} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

(3) Weight vector for 2nd layer $\tilde{W}^{(2)} = \begin{bmatrix} -1 \\ -1 \\ -1 \\ -1 \end{bmatrix}$

(4) Bias for 2nd layer, $b^{(2)} = 0$

Note: There can be many possible solutions of this problem!

Marking scheme: → A) Works for $x_1 > x_2 > x_3 > x_4 > x_5 \rightarrow 10/20$
 If your model → B) Works for $x_1 > x_2 > x_3 > x_4 > x_5 \rightarrow 20/20$

Problem 2: [20 pts] Consider a linear combination of M input variables x_1, \dots, x_D in the form

$$z_j = \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)}$$

where $j = 1, \dots, M$, resulting in the pre-activation z_j and the superscript (1) indicates that the corresponding parameters are in the first ‘layer’ of the network. Also consider

$$h_j = g(z_j)$$

where $g(\cdot)$ represents some differentiable, nonlinear activation function. The outputs of the activation functions are linearly combined to obtain the output pre-activation

$$a_k = \sum_{j=1}^M w_{kj}^{(2)} h_j + w_{k0}^{(2)}$$

where $k = 1, \dots, K$, and k is the total number of outputs. This transformation corresponds to the second layer of the network. Finally, the output unit activations are transformed using an appropriate activation function to give a set of network outputs y_k . The choice of activation function is determined by the nature of the data and the assumed distribution of target variables. Thus, for standard regression problems, the activation function is the identity so that $y_k = a_k$. Similarly, for multiple binary classification problems, each output unit activation is transformed using a logistic sigmoid function so that

$$y_k = \sigma(a_k)$$

where,

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

Finally, for multiclass problems, a SoftMax activation function is used. We can combine these various stages to give the overall network function that, for sigmoidal (logistic) output unit activation functions, takes the form

$$y_k(x; w) = \sigma \left(\sum_{j=1}^M w_{kj}^{(2)} g \left(\sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right)$$

where the set of all weight and bias parameters have been grouped together into a vector w . Thus, the neural network model is simply a nonlinear function from a set of input variables x_i to a set of output variables y_k controlled by a vector w of adjustable parameters. Considering a two-layer neural network as shown above with logistic sigmoid activation function, show that there exists an equivalent network,

which computes the same function, but with hidden unit activation functions given by $\tanh(a)$ where the $\tanh(a)$ is defined as follows:

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

Soln: Before showing how two networks (one with sigmoid activation at hidden layer, and another with tanh activation at hidden layer) are similar, establish the relation between sigmoid & tanh

Relation betw sigmoid & tanh:

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

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

$$\Rightarrow 2\sigma(x) = \tanh\left(\frac{x}{2}\right) + 1$$

$$\Rightarrow \boxed{\sigma(x) = \frac{1}{2} \tanh\left(\frac{x}{2}\right) + \frac{1}{2}}$$

Now let's look at the two networks

$$y_k(\tilde{x}; \tilde{w}) = \sigma \left(\sum_j w_{kj}^{(2)} g \left(\underbrace{\sum_i w_{ji}^{(1)} x_i + w_{j0}^{(1)}}_{\tilde{z}_j} \right) + w_{ko}^{(2)} \right)$$

For sigmoid activation $g = \sigma(\cdot)$

$$y_k(\tilde{x}; \tilde{w}) = \sigma \left(\sum_j w_{kj}^{(2)} \sigma(\tilde{z}_j) + w_{ko}^{(2)} \right)$$

For tanh activation, $g = \tanh(\cdot)$ [substitute the relation of $\sigma(\cdot)$ & $\tanh(\cdot)$]

$$y_k(\tilde{x}; \tilde{w}) = \sigma \left(\sum_j w_{kj}^{(2)} \left(\frac{1}{2} \tanh \left(\frac{\tilde{z}_j}{2} \right) + \frac{1}{2} \right) + w_{ko}^{(2)} \right)$$

$$= \sigma \left(\sum_j \frac{w_{kj}^{(2)}}{2} \tanh \left(\frac{\tilde{z}_j}{2} \right) + \sum_j \frac{w_{kj}^{(2)}}{2} + w_{ko}^{(2)} \right)$$

$$= \sigma \left(\sum_j \frac{w_{kj}^{(2)}}{2} \tanh \left(\sum_i \underbrace{w_{ji}^{(1)}}_{\frac{1}{2}} x_i + \underbrace{w_{j0}^{(1)}}_{\frac{1}{2}} \right) + \sum_j \underbrace{\frac{w_{kj}^{(2)}}{2} + w_{ko}^{(2)}}_{w_{kj}^{(2)*}} \right)$$

Parameters of equivalent network

EQUIVALENT Network with tanh hidden activation will have

$$y_k(\tilde{x}; \tilde{w}) = \sigma \left(\sum_j w_{kj}^{(2)*} \tanh \left(\sum_i w_{ji}^{(1)*} x_i + w_{j0}^{(1)*} \right) + w_{ko}^{(2)*} \right)$$

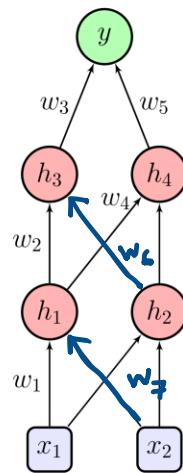
where,

$$w_{kj}^{(2)*} = \frac{w_{kj}^{(2)}}{2}, \quad w_{ko}^{(2)*} = w_{ko}^{(2)} + \sum_j \frac{w_{kj}^{(2)}}{2}, \quad w_{ji}^{(1)*} = \frac{w_{ji}^{(1)}}{2}, \quad w_{j0}^{(1)*} = \frac{w_{j0}^{(1)}}{2}$$

Marking scheme: 8 points for relation between $\sigma(\cdot)$ & $\tanh(\cdot)$

3 points each for $w_{ji}^{(1)*}$, $w_{j0}^{(1)*}$, $w_{kj}^{(2)*}$, $w_{ko}^{(2)*}$

Problem 3: [10 pts] One of the interesting features of the ReLU activation function is that it sparsifies the activations and the derivatives, i.e. sets a large fraction of the values to zero for any given input vector. Consider the following network:



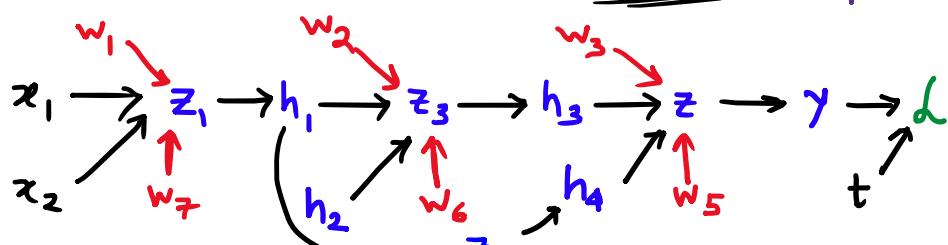
Note that each w_i refers to the weight on a single connection, not the whole layer. Suppose we are trying to minimize a loss function \mathcal{L} which depends only on the activation of the output unit y (For instance, \mathcal{L} could be the squared error loss $\frac{1}{2}(y - t)^2$). Suppose the unit h_1 receives an input of -1 on a particular training case, so the ReLU evaluates to 0 . Based only on this information, which of the weight derivatives

$$\frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial w_2}, \frac{\partial \mathcal{L}}{\partial w_3}$$

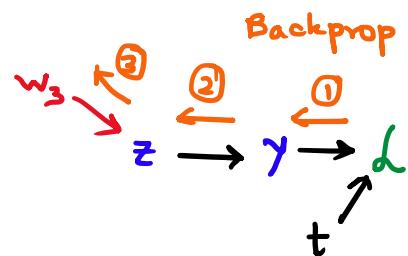
are guaranteed to be 0 for this training case? Write YES or NO for each. Justify your answers.

Soln: You may wish to define weights for two additional connections for convenience, as shown in figure (above).

Note: Each neuron is associated with a ReLU activation
So we could draw a partial computation graph (not necessary)



$$(a) \frac{\partial \mathcal{L}}{\partial w_3} = \underbrace{\frac{\partial \mathcal{L}}{\partial y}}_{\textcircled{1}} \times \underbrace{\frac{\partial y}{\partial z}}_{\textcircled{2}} \times \underbrace{\frac{\partial z}{\partial w_3}}_{\textcircled{3}}$$



During forward computation, the value of $z_1 = -1$ and $h_1 = 0$. All other variables need not be zero.

- $\frac{\partial \lambda}{\partial y}$ need not be 0
- $\frac{\partial y}{\partial z} = \frac{\partial}{\partial z} \text{ReLU}(z) = \frac{\partial}{\partial z} \max(0, z)$ need not be 0
- $\frac{\partial z}{\partial w_3} = \frac{\partial}{\partial w_3} (w_3 h_3 + w_5 h_4) = h_3 \leftarrow$ need not be 0
 $\therefore \frac{\partial \lambda}{\partial w_3}$ is not guaranteed to be zero

(b) $\frac{\partial \lambda}{\partial w_2} = \frac{\partial \lambda}{\partial z} \times \frac{\partial z}{\partial h_3} \times \frac{\partial h_3}{\partial z_3} \times \frac{\partial z_3}{\partial w_2}$

```

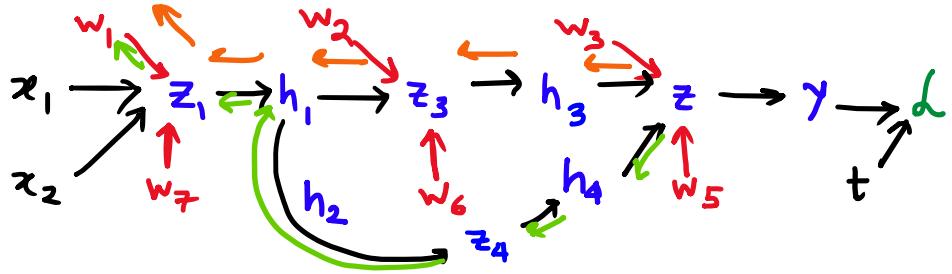
graph LR
    h1[h1] -- w2 --> z3[z3]
    h2[h2] -- w6 --> z3
    h3[h3] -- w3 --> z
    h4[h4] -- w5 --> z
    z --> y[y]
    y --> lambda[λ]
    t[t] --> lambda
  
```

$\frac{\partial \lambda}{\partial z} = \frac{\partial \lambda}{\partial y} \times \frac{\partial y}{\partial z}$ already proved to not necessarily zero

- $\frac{\partial z}{\partial h_3} = \frac{\partial}{\partial h_3} (w_3 h_3 + w_5 h_4) = w_3 \leftarrow$ need not be 0
- $\frac{\partial h_3}{\partial z_3} = \frac{\partial}{\partial z_3} \text{ReLU}(z_3) = \frac{\partial}{\partial z_3} \max(0, z_3) \leftarrow$ need not be 0
- $\frac{\partial z_3}{\partial w_2} = \frac{\partial}{\partial w_2} (w_2 h_1 + w_6 h_2) = h_1 = 0$ (given)

\therefore The product $\frac{\partial \lambda}{\partial w_2} = 0$ is guaranteed to be zero

(c)



$$\begin{aligned}
 \frac{\partial L}{\partial w_1} &= \underbrace{\frac{\partial L}{\partial z} \times \frac{\partial z}{\partial h_3} \times \frac{\partial h_3}{\partial z_3} \times \frac{\partial z_3}{\partial h_1}}_{A} \times \frac{\partial h_1}{\partial z_1} \times \frac{\partial z_1}{\partial w_1} \\
 &\quad + \underbrace{\frac{\partial L}{\partial z} \times \frac{\partial z}{\partial h_4} \times \frac{\partial h_4}{\partial z_4} \times \frac{\partial z_4}{\partial h_1}}_{B_1} \times \frac{\partial h_1}{\partial z_1} \times \frac{\partial z_1}{\partial w_1} \\
 &= \underbrace{\frac{\partial L}{\partial z}}_{B_2} \times \left(\underbrace{\frac{\partial z}{\partial h_3} \times \frac{\partial h_3}{\partial z_3} \times \frac{\partial z_3}{\partial h_1}}_{B_1} + \underbrace{\frac{\partial z}{\partial h_4} \times \frac{\partial h_4}{\partial z_4} \times \frac{\partial z_4}{\partial h_1}}_{B_2} \right) \times \underbrace{\frac{\partial h_1}{\partial z_1}}_{C} \times \underbrace{\frac{\partial z_1}{\partial w_1}}_{D}
 \end{aligned}$$

(A) $\frac{\partial L}{\partial z}$ ← need not be zero (proved previously)

- (B₁) • $\frac{\partial z}{\partial h_3} = \frac{\partial}{\partial h_3}(w_3 h_3 + w_5 h_5) = w_3$ ← need not be 0
- $\frac{\partial h_3}{\partial z_3} = \frac{\partial}{\partial z_3} \max(0, z_3)$ ← need not be 0
- $\frac{\partial z_3}{\partial h_1} = \frac{\partial}{\partial h_1} (w_2 h_1 + w_3 h_3) = w_2$ ← need not be 0

(B₂) Similarly you can check that $\frac{\partial z}{\partial h_4}$, $\frac{\partial h_4}{\partial z_4}$, $\frac{\partial z_4}{\partial h_1}$ need not be 0

(C) $\frac{\partial h_1}{\partial z_1} = \frac{\partial}{\partial z_1} \max(0, z_1) = \frac{\partial}{\partial z_1} \max(0, -1) = \frac{\partial}{\partial z_1} 0 = 0$

So the product turns out to be 0 because of $\frac{\partial h_1}{\partial z_1}$ being 0

$\therefore \frac{\partial L}{\partial w_1}$ is guaranteed to be 0

Problem 4: [20 pts] Consider a MLP consisting of N input units, K hidden units and N output units. The activations are computed as follows:

$$\begin{aligned} \underline{z} &= \mathbf{W}^{(1)} \underline{x} + \mathbf{b}^{(1)} \\ \underline{h} &= \sigma(\underline{z}) \\ \underline{y} &= \underline{x} + \mathbf{W}^{(2)} \underline{h} + \mathbf{b}^{(2)} \end{aligned}$$

where $\sigma(\cdot)$ is the sigmoid function, applied element wise. The loss \mathcal{L} will involve both \underline{h} and \underline{y} :

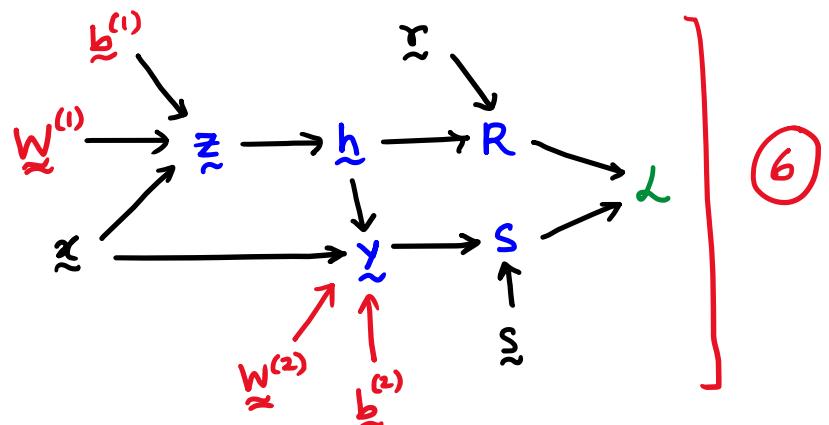
$$\begin{aligned} \mathcal{L} &= \mathcal{S} + \mathcal{R} \\ \mathcal{S} &= \frac{1}{2} \|\underline{y} - \underline{s}\|_2^2 = \frac{1}{2} (\underline{y} - \underline{s})^\top (\underline{y} - \underline{s}) \\ \mathcal{R} &= \underline{r}^T \underline{h} \end{aligned}$$

where \underline{r} and \underline{s} are given vectors.

1. Draw the computation graph relating \underline{x} , \underline{z} , \underline{h} , \underline{y} , \mathcal{R} , \mathcal{S} and \mathcal{L} .
2. Derive the backpropagation equations for computing $\bar{\underline{x}} = \frac{\partial \mathcal{L}}{\partial \underline{x}}$. You may use σ' to denote the derivation of the sigmoid function, so you do not have to write it out explicitly.

Solu:

1) Computation graph \rightarrow



$$2) \bar{\underline{x}} = \frac{\partial \mathcal{L}}{\partial \underline{x}} = 1, \quad \bar{R} = \frac{\partial \mathcal{L}}{\partial R} = 1, \quad \bar{S} = \frac{\partial \mathcal{L}}{\partial S} = 1$$

]

$$\bar{\underline{y}} = \frac{\partial \mathcal{L}}{\partial \underline{y}} = \bar{S} \frac{\partial S}{\partial \underline{y}} = \bar{S} \begin{bmatrix} \frac{\partial S}{\partial y_1} \\ \vdots \\ \frac{\partial S}{\partial y_n} \end{bmatrix} = \bar{S} \begin{bmatrix} y_1 - s_1 \\ \vdots \\ y_n - s_n \end{bmatrix} = \bar{S} \cdot (\underline{y} - \underline{s})$$

]

$$\begin{aligned}
\bar{h} &= \frac{\partial L}{\partial \bar{h}} = \frac{\partial R}{\partial \bar{h}} \bar{R} + \left(\frac{\partial \bar{y}}{\partial \bar{h}} \right)^T \bar{y} \\
&= \begin{bmatrix} \frac{\partial R}{\partial h_1} \\ \vdots \\ \frac{\partial R}{\partial h_K} \end{bmatrix} \bar{R} + \begin{bmatrix} w_{11}^{(2)} & \cdots & w_{1K}^{(2)} \\ \vdots & \ddots & \vdots \\ w_{n1}^{(2)} & \cdots & w_{nK}^{(2)} \end{bmatrix}^T \bar{y} \\
&= \begin{bmatrix} r_1 \\ \vdots \\ r_K \end{bmatrix} \bar{R} + \begin{bmatrix} w_{11}^{(2)} & \cdots & w_{1K}^{(2)} \\ \vdots & \ddots & \vdots \\ w_{n1}^{(2)} & \cdots & w_{nK}^{(2)} \end{bmatrix}^T \bar{y} \\
&= \bar{x} \bar{R} + \bar{w}^{(2)T} \bar{y} \quad] \textcircled{4}
\end{aligned}$$

$$\begin{aligned}
\bar{z} &= \frac{\partial L}{\partial z} = \left(\frac{\partial h}{\partial z} \right)^T \bar{h} = \begin{bmatrix} \frac{\partial h_1}{\partial z_1} & \cdots & \frac{\partial h_K}{\partial z_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_K}{\partial z_1} & \cdots & \frac{\partial h_K}{\partial z_K} \end{bmatrix}^T \bar{h} = \begin{bmatrix} \sigma'(z_1) & \cdots & \sigma'(z_K) \\ \vdots & \ddots & \vdots \\ \sigma'(z_1) & \cdots & \sigma'(z_K) \end{bmatrix}^T \bar{h} \\
&= \bar{h} \circ \sigma'(\bar{z}) \quad] \textcircled{2}
\end{aligned}$$

elementwise product

$$\begin{aligned}
\bar{x} &= \frac{\partial L}{\partial \bar{x}} = \left(\frac{\partial \bar{x}}{\partial \bar{x}} \right)^T \bar{z} + \left(\frac{\partial \bar{y}}{\partial \bar{x}} \right)^T \bar{y} \\
&= \bar{w}^{(1)T} \bar{z} + \bar{I} \bar{y} \\
&= \bar{w}^{(1)T} \bar{z} + \bar{y} \quad] \textcircled{3}
\end{aligned}$$