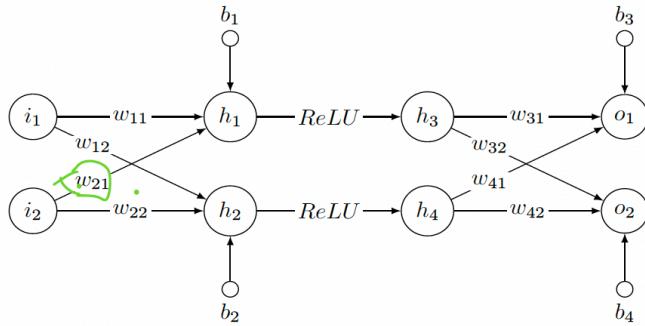


TUM

Part III: Backpropagation (9 points)

1. (9 points) Given the following neural network with fully connection layer and ReLU activations, including two input units (i_1, i_2), four hidden units (h_1, h_2) and (h_3, h_4). The output units are indicated as (o_1, o_2) and their targets are indicated as (t_1, t_2). The weights and bias of fully connected layer are called w and b with specific sub-descriptors.



The values of variables are given in the following table:

Variable	i_1	i_2	w_{11}	w_{12}	w_{21}	w_{22}	w_{31}	w_{32}	w_{41}	w_{42}	b_1	b_2	b_3	b_4	t_1	t_2
Value	2.0	-1.0	1.0	-0.5	0.5	-1.0	0.5	-1.0	-0.5	1.0	0.5	-0.5	-1.0	0.5	1.0	0.5

$$\begin{aligned}
 * o_1 &= \underbrace{\left(\underbrace{i_1 w_{11} + i_2 w_{21} + b_1}_{h_1} \right) \text{ReLU} \times w_{31}}_{\underbrace{h_3}_{h_2}} + \underbrace{\left(\underbrace{i_2 w_{22} + b_2}_{h_4} \right) \text{ReLU} \times w_{41}}_{\underbrace{h_4}_{h_2}} + b_3 \\
 &= \left((2.0 + (-1.0) \cdot 0.5 + 0.5) \text{ReLU} \times 0.5 \right) + (2.0 - 0.5) \\
 &\quad + (-1.0) (-1.0) + (-0.5) \text{ReLU} \times (0.5) + (-1.0) \\
 &= (2.0 - 0.5 + 0.5) \text{ReLU} \times 0.5 + (-1.0 + 1.0 - 0.5) \text{ReLU} \times -0.5 \\
 &= 1.0 - 1.0 = 0
 \end{aligned}$$

$$\begin{aligned}
 O_2 &= (i_1 w_{12} + i_2 w_{22} + b_2) \text{Relu} \times w_{42} + \text{Relu} \\
 &\quad ((i_1 w_{11} + i_2 w_{21} + b_1) \times w_{32} + b_3) \\
 &= -1.5
 \end{aligned}$$

$$\begin{aligned}
 L(O_1, O_2) &= \frac{1}{2} (O_1 - t_1)^2 + \frac{1}{2} (O_2 - t_2)^2 \\
 &= \frac{1}{2} (0 - 1)^2 + \frac{1}{2} (-1.5 - 0.5)^2 \\
 &= \frac{1}{2} (1 + 4) = 2.5
 \end{aligned}$$

* Update the weight w_{21} using gradient descent with learning rate 0.1.

$$\begin{aligned}
 \frac{\partial L}{\partial w_{21}} &= \frac{\partial L}{\partial o_1} \cdot \frac{\partial o_1}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_{21}} + \frac{\partial L}{\partial o_2} \cdot \frac{\partial o_2}{\partial h_3} \\
 &\quad \cdot \frac{\partial h_3}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_{21}} \\
 &= (o_1 - t_1) \cdot \frac{\partial}{\partial h_3} (h_3 w_{31} + h_1 w_{41} + b_3)
 \end{aligned}$$

$$\begin{aligned}
& \cdot \frac{\partial}{\partial h_1} (\text{ReLU}(h_1)) \cdot \frac{\partial}{\partial w_{21}} (i_1 w_{11} + i_2 w_{21} - b_1) \\
& + (o_2 - t_2) \cdot \frac{\partial}{\partial h_3} (h_3 w_{32} + h_4 w_{42} + b_4) \\
& \cdot \frac{\partial}{\partial h_1} (\text{ReLU}(h_1)) \cdot \frac{\partial}{\partial w_{21}} (i_1 w_{11} + i_2 w_{21} + b_1) \\
= & 0.5 \cdot w_{31} \cdot \text{ReLU}'(h_1) \cdot i_2 + (o_2 - t_2) \\
& \cdot w_{32} \cdot \text{ReLU}'(h_1) \cdot i_2
\end{aligned}$$

$$= -1 \times 0.5 \times 1 \cdot (-1) + (-2) \cdot (-1.0) \cdot 1 \cdot (-1.0)$$

$$= 0.5 - 2$$

$$= -1.5$$

$$\begin{aligned}
w_{21}^{\text{new}} &= w_{21} \text{old} - n \nabla \frac{dL}{dw_{21}} \\
&= 0.5 - 0.1 \cdot (-1.5)
\end{aligned}$$

$$= \frac{5}{10} + \frac{15}{100}$$

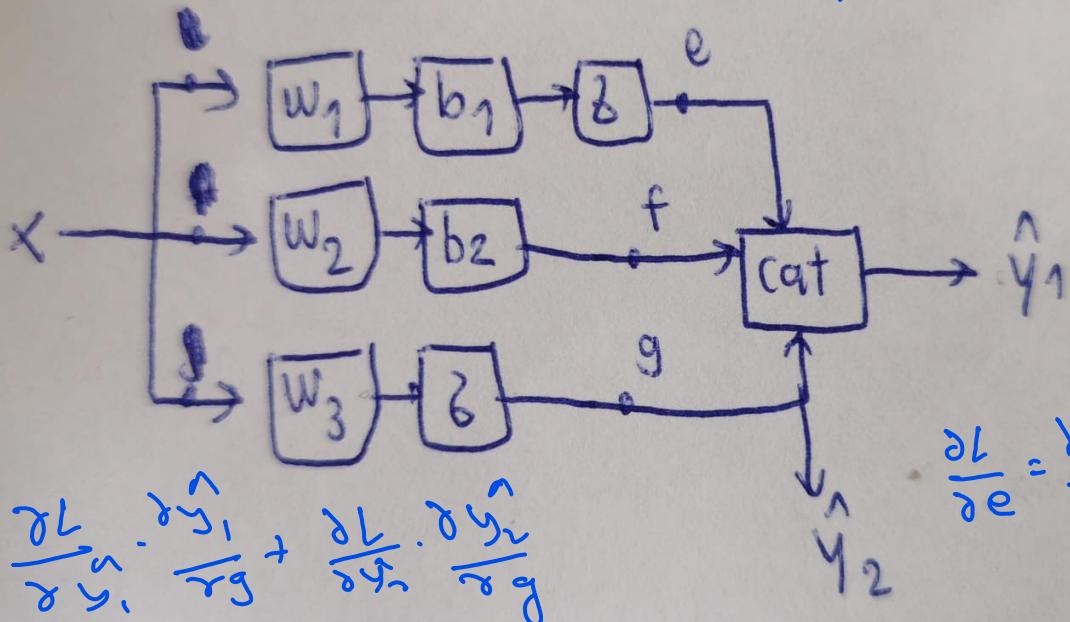
$$= \frac{50 + 15}{100} = 0.65$$

53-23

$$\frac{\partial L}{\partial \hat{y}_1} = d(\hat{y}_1 - y)^2 = 2(\hat{y}_1 - y)$$

$$\frac{\partial L}{\partial \hat{y}_2} = (\hat{y}_2 - y)^2 = 2(\hat{y}_2 - y)$$

$$MSE = (\hat{y} - y)^2$$



$$e = g(w_1 x + b_1)$$

$$e = (x w_1 + b_1) g$$

$$f = x w_2 + b_2$$

$$g = (x w_3) g$$

$$f = \cancel{x w_2} + b_2$$

$$\hat{y}_2 = g = g w_3 x$$

$$\hat{y}_1 = e + f + g$$

$$\hat{y}_1 = e + f + g$$

$$\hat{y}_2 = g$$

$$\frac{\partial L}{\partial \hat{y}_1} = \frac{\partial}{\partial \hat{y}_1} (\hat{y}_1 - y)^2$$

$$= 2(\hat{y}_1 - y)$$

$$\frac{\partial L}{\partial \hat{y}_2} = \frac{\partial}{\partial \hat{y}_2} (\hat{y}_2 - y)^2$$

$$= 2(\hat{y}_2 - y)$$

$$\frac{\partial L}{\partial e} = \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial e}$$

$$= \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial}{\partial e} (e + f + g)$$

$$= 2(\hat{y}_1 - y) \cdot 1$$

$$\frac{\partial L}{\partial f} = \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial f}$$

$$= \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial}{\partial f} (e + f + g)$$

$$= 2(\hat{y}_1 - y) \cdot 1$$

$$\frac{\partial L}{\partial g} = \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial g} + \frac{\partial L}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial g}$$

$$= \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial}{\partial g} (e + f + g) + \frac{\partial L}{\partial \hat{y}_2} \cdot \frac{\partial}{\partial g} (g)$$

$$= 2(\hat{y}_1 - y) \cdot 1 + 2(\hat{y}_2 - y) \cdot 1$$

$$\begin{aligned}\frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial e} \cdot \frac{\partial e}{\partial w_1} \\ &= \frac{\partial L}{\partial e} \cdot \frac{\partial}{\partial w_1} \left(g(xw_1 + b_1) \right) \\ &= \frac{\partial L}{\partial e} \cdot g'(xw_1 + b_1) \cdot x\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial w_2} &= \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial w_2} \\ &= \frac{\partial L}{\partial f} \cdot \frac{\partial}{\partial w_2} \left(g(xw_2 + b_2) \right) \\ &= \frac{\partial L}{\partial f} \cdot x\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial w_3} &= \frac{\partial L}{\partial g} \cdot \frac{\partial g}{\partial w_3} \\ &= \frac{\partial L}{\partial g} \cdot \frac{\partial}{\partial w_3} \left(g(xw_3) \right) \\ &= \frac{\partial L}{\partial g} \cdot g'(xw_3) \cdot x\end{aligned}$$

Mock - 23 July 2021

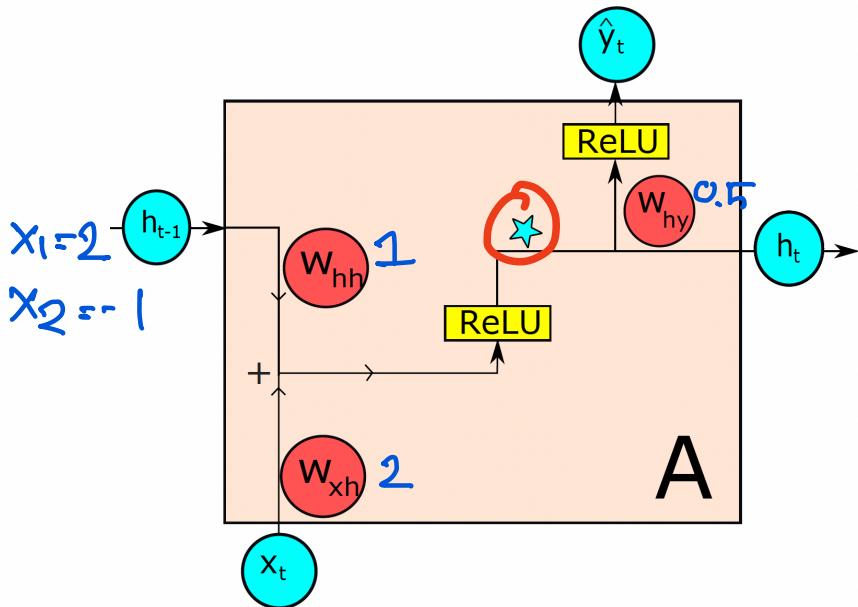


Figure 3: Schematic of a simplified Elman Cell.

* Calculate $h_1, h_2, \hat{y}_1, \hat{y}_2$ for $t=1$ and 2 .

$$h_1 = \text{ReLU}(h_{t-1} w_{hh} + w_{xh} x_t)$$

$$= \text{ReLU}(0 \cdot w_{hh} + 2 \cdot 2)$$

$$= 4$$

$t=1$

$$h_2 = \text{ReLU}(h_{t-1} w_{hh} + x_t w_{xh})$$

$t=2$

$$= \text{ReLU}(4 \cdot 1 + (-1) \cdot 2)$$

$$= \text{ReLU}(2)$$

$$= 2$$

$$\begin{aligned}\hat{y}_1 &= \text{ReLU}(h_1 \times w_{hy}) \\ &= \text{ReLU}(c_1 \times 0.5) \\ &= 1\end{aligned}$$

$$\begin{aligned}\hat{y}_2 &= \text{ReLU}(h_2 \times w_{hy}) \\ &= \text{ReLU}(2 \times 0.5) \\ &= 1\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial h_t} &= \frac{\partial L}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial h_t} + \frac{\partial L}{\partial h_{t+1}} \cdot \frac{\partial h_{t+1}}{\partial h_t} \\ &= 2(\hat{y}_t - y) \cdot \frac{\partial}{\partial h_t} \left(\text{ReLU}(h_t \times w_{hy}) \right) \\ &\quad + \frac{\partial L}{\partial h_{t+1}} \cdot \frac{\partial h_{t+1}}{\partial h_t} \\ &= \frac{\partial L}{\partial \hat{y}_t} \cdot \text{ReLU}'(w_{hy}, h_t) \cdot w_{hy} \\ &\quad + \frac{\partial L}{\partial h_{t+1}} \cdot \frac{\partial}{\partial h_t} \left(\text{ReLU}(h_t \cdot w_{hn} + b_{t+1} \cdot w_{xh}) \right)\end{aligned}$$

$$= \frac{\partial L}{\partial g^+} \cdot \text{Relu}'(w_{hy}, h^+) \cdot w_{hy}$$

$$+ \frac{\partial L}{\partial h_{t+1}} \cdot \text{Relu}'(h^+_t \cdot w_{hh} + x^+_t \cdot w_{xh}) \cdot \underbrace{w_{hh}}_{y_{t+1}}$$

$$\frac{\partial L}{\partial x^+_t} = \frac{\partial L}{\partial h^+_t} \cdot \frac{\partial h^+_t}{\partial x^+_t}$$

$$= \frac{\partial L}{\partial h^+_t} \cdot \frac{\partial}{\partial x^+_t} \left(\text{Relu}(h^+_{t+1} \cdot w_{hh} + x^+_t \cdot w_{xh}) \right)$$

$$= \frac{\partial L}{\partial h^+_t} \cdot \text{Relu}'(h^+_{t+1} \cdot w_{hh} + x^+_t \cdot w_{xh}) \cdot w_{xh}$$

$$\frac{\partial L}{\partial w_{xh}} = \frac{\partial L}{\partial h^+_t} \cdot \frac{\partial h^+_t}{\partial w_{xh}}$$

$$= \frac{\partial L}{\partial h^+_t} \cdot \frac{\partial}{\partial w_{xh}} \left(\text{Relu}(h^+_{t+1} \cdot w_{hh} + x^+_t \cdot w_{xh}) \right)$$

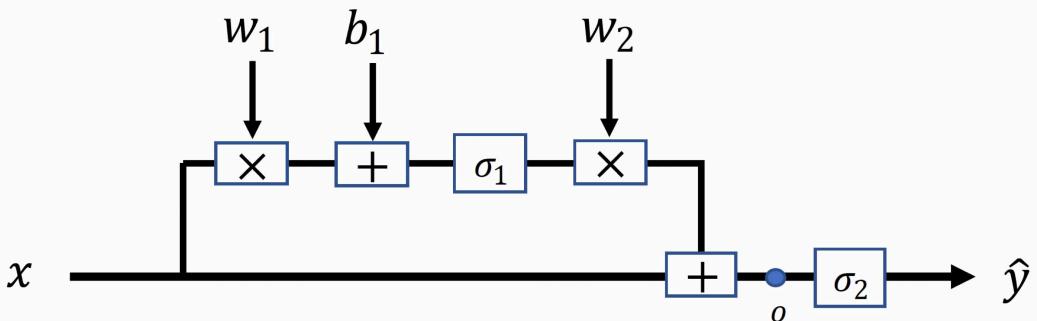
$$= \frac{\partial L}{\partial h^+_t} \cdot \text{Relu}'(h^+_{t+1} \cdot w_{hh} + x^+_t \cdot w_{xh}) \cdot x^+_t$$

$$\frac{\partial L}{\partial w_{hy}} = \frac{\partial L}{\partial h^+_t} \cdot \frac{\partial h^+_t}{\partial w_{hy}} = \frac{\partial L}{\partial h^+_t} \cdot \frac{\partial}{\partial w_{hy}} \left(\text{Relu}(h^+_{t+1} \cdot w_{hh} + x^+_t \cdot w_{xh}) \right)$$

$$= \frac{\partial L}{\partial h^+_t} \cdot \text{Relu}'(h^+_{t+1} \cdot w_{hy} + x^+_t \cdot w_{xh}) \cdot h^+_{t+1}$$

26 July 2021

Given is the following network $f(x) = \hat{y}$ receiving $x \in \mathbb{R}$ as input to compute a prediction $\hat{y} \in \mathbb{R}$. It uses two weights $w_1, w_2 \in \mathbb{R}$, one bias $b_1 \in \mathbb{R}$ and two sigmoid activations denoted as functions $\sigma_1(\cdot)$ and $\sigma_2(\cdot)$ shown in the figure below. The network is trained using the L_2 norm, defined as $L(y, \hat{y}) = \|\hat{y} - y\|_2^2$ with labels $y \in \{0, 1\}$. The boxes are mathematical operations, while \times represents the multiplication, $+$ the addition and σ the sigmoid activation. o marks an intermediate result as indicated in the figure.



$$\begin{aligned}
 f(x) &= \hat{y} = ((w_1 x + b_1) \sigma_1 w_2 + x) \sigma_2 \\
 &= \left\{ \underbrace{0.5 \cdot 1 + (-0.5)}_{=} \right\} \sigma_1 \left\{ \underbrace{x \cdot (-2) + 1}_{=} \right\} \sigma_2 \\
 &= 0.5
 \end{aligned}$$

$$\begin{aligned}
 L &= (\hat{y} - y)^2 \\
 &= (0.5 - 1)^2 \\
 &= 0.25
 \end{aligned}$$

$$\frac{\partial L}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}} (\hat{y} - y)^2$$

$$= 2(\hat{y} - y)$$



$$\frac{\partial L}{\partial o} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o}$$

$$= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial}{\partial o} (o_2(o))$$

$$= \frac{\partial L}{\partial \hat{y}} \cdot o_2'(o)$$



$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial o} \cdot \frac{\partial o}{\partial w_2}$$

$$= \frac{\partial L}{\partial o} \cdot \frac{\partial}{\partial w_2} (x \cdot \{ (w_1 \cdot x + b_1) \cdot o_1, x \cdot w_2 \})$$

$$= \frac{\partial L}{\partial o} \cdot o_1(w_1 \cdot x + b_1)$$



$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial o} \cdot \frac{\partial o}{\partial b_1}$$

$$= \frac{\partial L}{\partial o} \cdot \frac{\partial}{\partial b_1} (x \cdot \{ (w_1 \cdot x + b_1) \cdot o_1, x \cdot w_2 \})$$

$$= \frac{\partial L}{\partial o} \cdot o_1'(w_1 \cdot x + b_1)$$



$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial o} \cdot \frac{\partial o}{\partial w_1}$$

$$= \frac{\partial L}{\partial o} \cdot \frac{\partial}{\partial w_1} \left(x_{+}(\{w_1x + b_1\}g_1, \{w_2\}) \right)$$

$$= \frac{\partial L}{\partial o} \cdot w_2 \cdot g_1'(\{w_1x + b_1\}) \cdot x$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial o} \cdot \frac{\partial o}{\partial x}$$

$$= \frac{\partial L}{\partial o} \cdot \frac{\partial}{\partial x} \left(x_{+}(\{w_1x + b_1\}g_1, \{w_2\}) \right)$$

$$= \frac{\partial L}{\partial o} \cdot (1 + w_2 \cdot g_1'(\{w_1x + b_1\}) \cdot w_1)$$

$$= \frac{\partial L}{\partial o} + \frac{\partial L}{\partial o} \cdot w_2 \cdot g_1'(\{w_1x + b_1\}) w_1$$

24 June 2021

4 Backpropagation & Recurrent Neural Networks (9P)

Backpropagation is an important concept that enables training neural networks and it is also used in recurrent neural networks. Given is the recurrent cell visualized in Figure 2. In the following, the subscript (e.g., x_t) will denote the time step t .

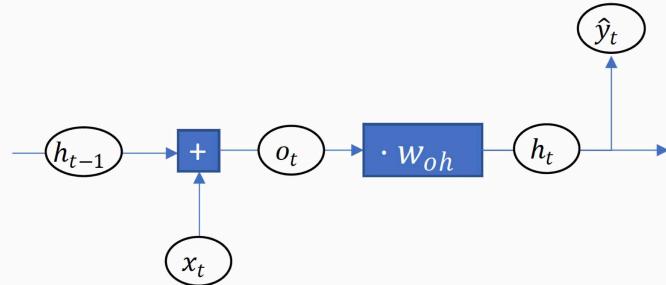


Figure 2: Schematic of a novel recurrent cell.

$$h_1 = (h_{t-1} + x_t) \omega_{oh}$$

$$= (0 + 1) 0.5$$

$$= 0.5$$

$$h_2 = (h_{t-1} + x_t) \omega_{oh}$$

$$= (0.5 + 0.5) 0.5$$

$$= 0.5$$

$$\hat{y}_1 = h_t = 0.5$$

$$\hat{y}_2 = h_t = 0.5$$

$$L = -y_+ \ln(\hat{y}_+) - (1-y_+) \ln(1-\hat{y}_+)$$

$$\frac{\partial L}{\partial \hat{y}_+} = -\frac{y}{\hat{y}_+} + \frac{1-y}{1-\hat{y}_+}$$

$$\begin{aligned}\frac{\partial L}{\partial h_t} &= \frac{\partial L}{\partial \hat{y}_+} \cdot \frac{\partial \hat{y}_+}{\partial h_t} + \frac{\partial L}{\partial h_{t+1}} \cdot \frac{\partial h_{t+1}}{\partial h_t} \\ &= \frac{\partial L}{\partial \hat{y}_+} \cdot \frac{\partial}{\partial h_t} (h_+) + \frac{\partial L}{\partial h_{t+1}} \cdot \frac{\partial}{\partial h_t} (h_{t+1} + x_t) \\ &= \frac{\partial L}{\partial \hat{y}_+} \cdot 1 + \frac{\partial L}{\partial h_{t+1}} \cdot w_{oh}\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial o_t} &= \frac{\partial L}{\partial h_t} \cdot \frac{\partial h_t}{\partial o_t} \\ &= \frac{\partial L}{\partial h_t} \cdot \frac{\partial}{\partial o_t} (o_t \times w_{on}) \\ &= \frac{\partial L}{\partial h_t} \cdot w_{on}\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial w_{oh}} &= \frac{\partial L}{\partial h_t} \cdot \frac{\partial h_t}{\partial w_{oh}} \\ &= \frac{\partial L}{\partial h_t} \cdot \frac{\partial}{\partial w_{oh}} (o_t \times w_{on})\end{aligned}$$

$$= \frac{\partial L}{\partial h_t} \cdot o_t$$

$$\frac{\partial L}{\partial W_{ohn}} = \sum_f \frac{\partial L}{\partial W_{ohn}}$$

29 March 2023

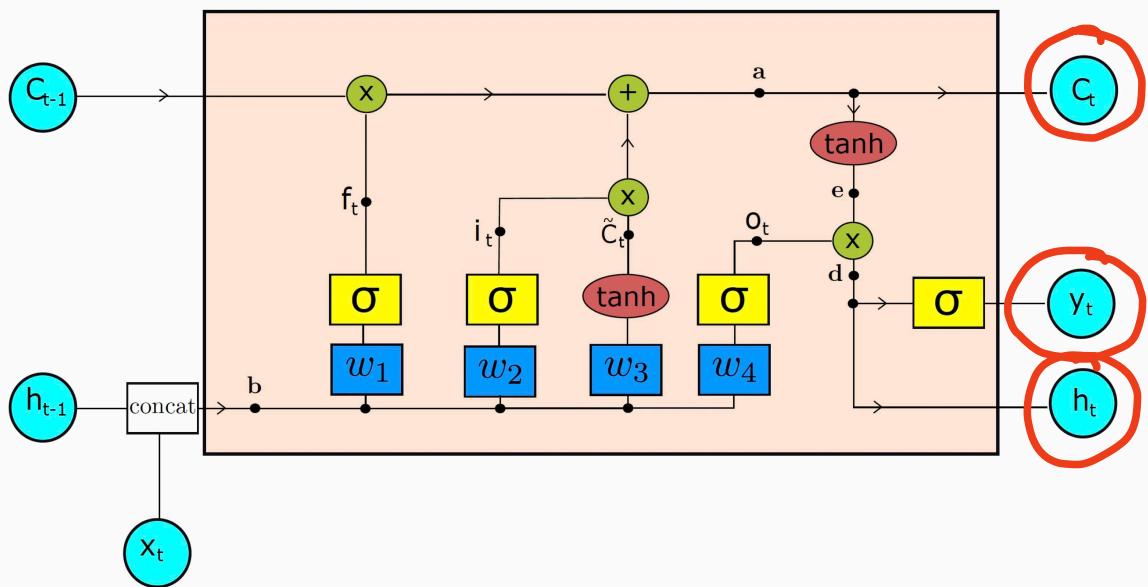


Figure 2: LSTM cell with intermediate steps.

C_t = long term memory
 h_t = short term memory
 \hat{y}_t = output

$$\tanh'(x) = 1 - \tanh^2(x)$$

$$G'(x) = G(x)(1-G(x))$$

$$C_t = (C_{t-1} \cdot f_t) + (\hat{o}_t \cdot \hat{c}_t)$$

$$Y_t = (\tanh(C_t) \cdot o_t)^6$$

$$h_t = \tanh(C_t) \cdot o_t$$

$$\frac{\partial L}{\partial C_{t-1}} = \frac{\partial L}{\partial C_t} \cdot \frac{\partial C_t}{\partial C_{t-1}}$$

$$= \frac{\partial L}{\partial C_t} \frac{\partial}{\partial C_{t-1}} \left((C_{t-1} \cdot f_t) + (\hat{o}_t \times \hat{c}_t) \right)$$

$$= \frac{\partial L}{\partial C_t} \cdot f_t$$

$$\frac{\partial L}{\partial d} = \frac{\partial L}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial d} + \frac{\partial L}{\partial h_t} \cdot \frac{\partial h_t}{\partial d}$$

$$= \frac{\partial L}{\partial y+} \cdot \frac{\partial}{\partial d} (6(d)) + \frac{\partial L}{\partial h+} \cdot \frac{\partial}{\partial d} (d)$$

$$= \frac{\partial L}{\partial y+} \cdot 6'(d) + \frac{\partial L}{\partial h+} \cdot 1$$

$$\frac{\partial L}{\partial e} = \frac{\partial L}{\partial d} \cdot \frac{\partial d}{\partial e}$$

$$= \frac{\partial L}{\partial d} \cdot \frac{\partial}{\partial e} (e \cdot o+)$$

$$= \frac{\partial L}{\partial d} \cdot o+$$

$$\frac{\partial L}{\partial o+} = \frac{\partial L}{\partial d} \cdot \frac{\partial d}{\partial o+}$$

$$= \frac{\partial L}{\partial d} \cdot \frac{\partial}{\partial o+} (o+ \cdot e)$$

$$= \frac{\partial L}{\partial d} \cdot e$$

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial e} \cdot \frac{\partial e}{\partial a} + \frac{\partial L}{\partial c+}$$

$$= \frac{\partial L}{\partial e} \cdot \frac{\partial}{\partial a} (\tanh(a)) + \frac{\partial L}{\partial c+}$$

$$= \frac{\partial L}{\partial e} \cdot \tanh'(a) + \frac{\partial L}{\partial c+}$$

$$\frac{\partial L}{\partial \hat{c}^+} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial \hat{c}^+}$$

$$= \frac{\partial L}{\partial a} \cdot \frac{\partial}{\partial \hat{c}^+} \left(C_{t-1} \cdot f_t + \hat{o}_t \cdot \hat{c}^+ \right)$$

$$= \frac{\partial L}{\partial a} \cdot \hat{o}_t$$

$$\frac{\partial L}{\partial \hat{o}_t} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial \hat{o}_t}$$

$$= \frac{\partial L}{\partial a} \cdot \frac{\partial}{\partial \hat{o}_t} \left(C_{t-1} \cdot f_t + \hat{o}_t \cdot \hat{c}^+ \right)$$

$$= \frac{\partial L}{\partial a} \cdot \hat{c}^+$$

$$\frac{\partial L}{\partial f_t} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial f_t} \left(C_{t-1} \cdot f_t + \hat{o}_t \cdot \hat{c}^+ \right)$$

$$= \frac{\partial L}{\partial a} \cdot C_{t-1}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial f_t} \cdot \frac{\partial f_t}{\partial b} + \frac{\partial L}{\partial \hat{o}_t} \cdot \frac{\partial \hat{o}_t}{\partial b} + \frac{\partial L}{\partial \hat{c}^+} \cdot \frac{\partial \hat{c}^+}{\partial b}$$

$$+ \frac{\partial L}{\partial u_t} \cdot \frac{\partial u_t}{\partial b}$$

$$= \frac{\partial L}{\partial f_f} \cdot \frac{\partial}{\partial b} (g(bw_1)) + \frac{\partial L}{\partial i_f} \cdot \frac{\partial}{\partial b} (g(bw_2))$$

$$+ \frac{\partial L}{\partial \tilde{e}_f} \cdot \frac{\partial}{\partial b} (\tanh(bw_3)) + \frac{\partial L}{\partial o_f} \cdot \frac{\partial}{\partial b} (g(bw_4))$$

$$= \frac{\partial L}{\partial f_f} \cdot g'(bw_1) \cdot w_1 + \frac{\partial L}{\partial i_f} \cdot g'(bw_2) \cdot w_2$$

$$+ \frac{\partial L}{\partial \tilde{e}_f} \cdot \tanh'(bw_3) \cdot w_3 + \frac{\partial L}{\partial o_f} \cdot g'(bw_4) \cdot w_4$$

08 August 2022

4 Recurrent Networks and Backpropagation (9P)

Given is the following network which receives an input $x \in \mathbb{R}$ to predict two outputs $f_1(x) = \hat{y}_1$ and $f_2(x) = \hat{y}_2$ where $\hat{y}_{1,2} \in \mathbb{R}$. It only contains one weight $w \in \mathbb{R}$ which is multiplied with its input and a bias $b \in \mathbb{R}$ which is added to its input. At the end a sigmoid function σ is applied. The states k , l and m are highlighted in the figure as intermediate results. The + operation for the skip connection is a simple addition. The network is visualized in the following figure:

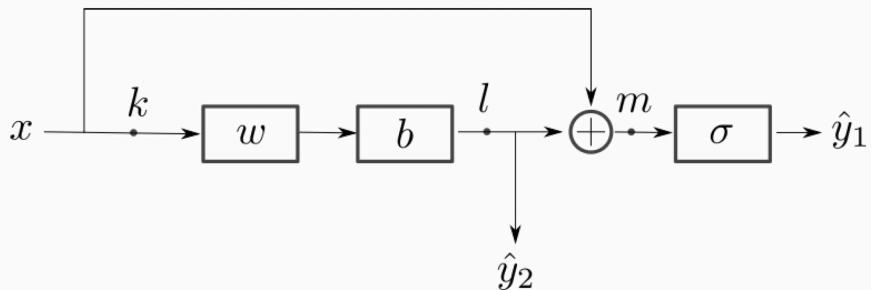


Figure 1: Schematic of a multi-headed residual network.

$$l = wx + b$$

$$m = wx + b + x$$

$$f_1(x) = \hat{y}_1 = \sigma(m) = \sigma(wx + b + x)$$

$$f_2(x) = \hat{y}_2 = \sigma(m)$$

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

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

$$\begin{aligned}
 \frac{\partial L}{\partial m} &= \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial m} \\
 &= \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial}{\partial m} (\sigma(m)) \\
 &= \frac{\partial L}{\partial \hat{y}_1} \cdot \sigma'(m)
 \end{aligned}$$

$$\frac{\partial L}{\partial l} = \frac{\partial L}{\partial m} \cdot \frac{\partial m}{\partial l} + \frac{\partial L}{\partial g_2} \cdot \frac{\partial g_2}{\partial l}$$

$$= \frac{\partial L}{\partial m} \cdot \frac{\partial}{\partial l} (l+x) + \frac{\partial L}{\partial g_2} \cdot \frac{\partial}{\partial l} (u)$$

$$= \frac{\partial L}{\partial m} \cdot 1 + \frac{\partial L}{\partial g_2} \cdot 1$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial l} \cdot \frac{\partial l}{\partial b}$$

$$= \frac{\partial L}{\partial l} \cdot \frac{\partial}{\partial b} (xw+b)$$

$$= \frac{\partial L}{\partial l} \cdot 1$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial l} \cdot \frac{\partial l}{\partial w}$$

$$= \frac{\partial L}{\partial l} \cdot \frac{\partial}{\partial w} (xw+b)$$

$$= \frac{\partial L}{\partial l} \cdot x$$

$$\frac{\partial L}{\partial k} = \frac{\partial L}{\partial l} \cdot \frac{\partial l}{\partial k}$$

$$= \frac{\partial L}{\partial l} \cdot \frac{\partial}{\partial k} (kxw+b)$$

$$= \frac{\partial L}{\partial e} \cdot w$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial k} \cdot \frac{\partial k}{\partial x} + \frac{\partial L}{\partial m} \cdot \frac{\partial m}{\partial x}$$

$$= \frac{\partial L}{\partial k} \cdot \frac{\partial}{\partial x}(x) + \frac{\partial L}{\partial m} \cdot \frac{\partial}{\partial x}(x)$$

$$= \frac{\partial L}{\partial k} \cdot 1 + \frac{\partial L}{\partial m} \cdot 1$$

[
m = x
just fun
short
memo
py]

06 April 2022

4 Recurrent Networks and Backpropagation (14 P)

Given is the following modified recurrent network cell $f(x_t, h_{t-1}) = \hat{y}_t$. It receives an input $x_t \in \mathbb{R}$ and a hidden state $h_t \in \mathbb{R}$ to compute a prediction \hat{y}_t . It only contains one weight $w \in \mathbb{R}$ which is multiplied with its input. The states s_t and o_t are highlighted in the figure as intermediate results. The $+$ operation is a simple addition. Using the label $y_t \in \{0, 1\}$ and the loss function $L(y_t, \hat{y}_t) = \|\hat{y}_t - y_t\|_2^2$ one can compute the loss and backpropagate it through the network for an input sequence of length $T \in \mathbb{N}^+$ with $t \in [1, T]$ being the current time state. The network is visualized in the following figure:

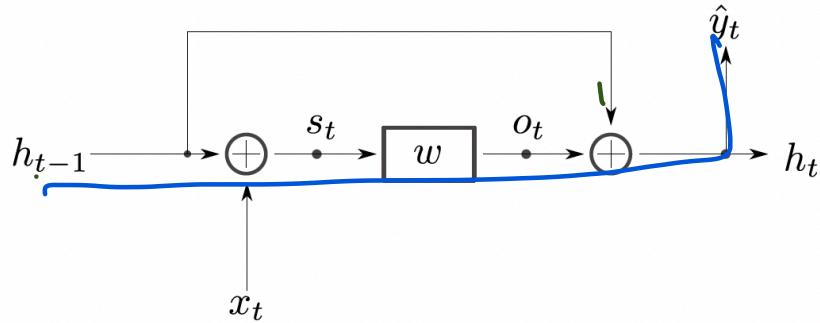


Figure 2: Schematic of a recurrent network cell.

$$s_t = h_{t-1} + x_t$$

$$o_t = (h_{t-1} + x_t) \omega$$

$$f = (h_{t-1} + x_t) \omega + h_{t-1}$$

$$\frac{\partial L}{\partial \hat{y}_t} = \frac{\partial}{\partial \hat{y}_t} (\hat{y}_t - y)^2$$

$$= 2(\hat{y}_t - y)$$

$$\frac{\partial L}{\partial o_t} = \frac{\partial L}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial o_t} + \frac{\partial L}{\partial h_t}$$

$$= \frac{\partial L}{\partial \hat{y}_t} \cdot \frac{\partial}{\partial o_t} (o_t + h_{t-1}) + \frac{\partial L}{\partial h_t}$$

$$= \frac{\partial L}{\partial y_+} \cdot 1 + \frac{\partial L}{\partial h_+}$$

$$\frac{\partial L}{\partial w_+} = \frac{\partial L}{\partial o_+} \cdot \frac{\partial o_+}{\partial w_+}$$

$$= \frac{\partial L}{\partial o_+} \cdot \frac{\partial}{\partial w_+} (st. w)$$

$$= \frac{\partial L}{\partial o_+} \cdot b_t$$

$$\frac{\partial L}{\partial w} = \sum \frac{\partial L}{\partial w_+}$$

$$\frac{\partial L}{\partial s_+} = \frac{\partial L}{\partial o_+} \cdot \frac{\partial o_+}{\partial s_+}$$

$$= \frac{\partial L}{\partial o_+} \cdot \frac{\partial}{\partial s_+} (st. w)$$

$$= \frac{\partial L}{\partial o_+} \cdot w$$

$$\frac{\partial L}{\partial x_+} = \frac{\partial L}{\partial s_+} \cdot \frac{\partial s_+}{\partial x_+}$$

$$= \frac{\partial L}{\partial s_+} \cdot \frac{\partial}{\partial x_+} (h_{+1} + x_+)$$

$$= \frac{\partial L}{\partial s_+} \cdot 1$$

$$\frac{\partial L}{\partial h_{+1}} = \frac{\partial L}{\partial \delta t} \cdot \frac{\partial \delta t}{\partial h_{+1}} + \frac{\partial L}{\partial h_t} \cdot \frac{\partial h_t}{\partial h_{+1}} + \frac{\partial L}{\partial \bar{y}_t} \cdot \frac{\partial \bar{y}_t}{\partial h_{+1}}$$

$$= \frac{\partial L}{\partial \delta t} \cdot \frac{\partial}{\partial h_{+1}} (h_{+1} + x_+) + \frac{\partial L}{\partial h_t} \cdot \frac{\partial}{\partial h_{+1}} ($$

$$(h_{+1} + o_t) + \frac{\partial L}{\partial \bar{y}_t} \cdot \frac{\partial}{\partial h_{+1}} (h_{+1} + o_t)$$

$$= \frac{\partial L}{\partial \delta t} \cdot 1 + \frac{\partial L}{\partial h_t} \cdot 1 + \frac{\partial L}{\partial \bar{y}_t} \cdot 1$$

$$\underline{\frac{\partial L}{\partial o_t}}$$

$$= \frac{\partial L}{\partial \delta t} + \frac{\partial L}{\partial o_t}$$