

1. a)

1. a) By Chain Rule

$$\begin{aligned} \bar{h}^{(t)} &= \bar{i}^{(t+1)} \frac{\partial i^{(t+1)}}{\partial h^{(t)}} + \bar{g}^{(t+1)} \frac{\partial g^{(t+1)}}{\partial h^{(t)}} + \bar{f}^{(t+1)} \frac{\partial f^{(t+1)}}{\partial h^{(t)}} + \bar{o}^{(t+1)} \frac{\partial o^{(t+1)}}{\partial h^{(t)}} \\ &= \bar{i}^{(t+1)} \frac{\partial i^{(t+1)}}{\partial w_{ix}^{(t+1)} + w_{ih}^{(t+1)}} \cdot \frac{\partial h^{(t+1)}}{\partial h^{(t)}} + \bar{g}^{(t+1)} \frac{\partial g^{(t+1)}}{\partial w_{gx}^{(t+1)} + w_{gh}^{(t+1)}} \cdot \frac{\partial h^{(t+1)}}{\partial h^{(t)}} \\ &\quad + \bar{f}^{(t+1)} \frac{\partial f^{(t+1)}}{\partial w_{fx}^{(t+1)} + w_{fh}^{(t+1)}} \cdot \frac{\partial h^{(t+1)}}{\partial h^{(t)}} + \bar{o}^{(t+1)} \frac{\partial o^{(t+1)}}{\partial w_{ox}^{(t+1)} + w_{oh}^{(t+1)}} \cdot \frac{\partial h^{(t+1)}}{\partial h^{(t)}} \\ &= \bar{i}^{(t+1)} \sigma'(w_{ix}^{(t+1)} + w_{ih}^{(t+1)}) \bar{h}^{(t)} + \bar{g}^{(t+1)} \sigma'(w_{gx}^{(t+1)} + w_{gh}^{(t+1)}) \bar{h}^{(t)} \\ &\quad + \bar{f}^{(t+1)} \sigma'(w_{fx}^{(t+1)} + w_{fh}^{(t+1)}) \bar{h}^{(t)} + \bar{o}^{(t+1)} \sigma'(w_{ox}^{(t+1)} + w_{oh}^{(t+1)}) \bar{h}^{(t)} \\ \bar{c}^{(t)} &= \bar{h}^{(t)} \frac{\partial h^{(t)}}{\partial \tanh(c^{(t)})} \cdot \frac{\partial \tanh(c^{(t)})}{\partial c^{(t)}} + \bar{c}^{(t+1)} \frac{\partial c^{(t+1)}}{\partial c^{(t)}} \\ &= \bar{h}^{(t)} \sigma^{(t)} \cdot \tanh'(c^{(t)}) + \bar{c}^{(t+1)} f^{(t+1)} \\ \bar{g}^{(t)} &= \bar{c}^{(t)} \frac{\partial c^{(t)}}{\partial g^{(t)}} \quad ; \quad \bar{o}^{(t)} = \bar{h}^{(t)} \frac{\partial h^{(t)}}{\partial o^{(t)}} \quad ; \quad \bar{f}^{(t)} = \bar{c}^{(t)} \frac{\partial c^{(t)}}{\partial f^{(t)}} \quad ; \quad \bar{i}^{(t)} = \bar{c}^{(t)} \frac{\partial c^{(t)}}{\partial i^{(t)}} \\ &= \bar{c}^{(t)} \cdot \bar{i}^{(t+1)} \quad ; \quad = \bar{h}^{(t)} \tanh(c^{(t)}) \quad ; \quad = \bar{c}^{(t)} \cdot \bar{c}^{(t+1)} \quad ; \quad = \bar{c}^{(t)} \cdot \bar{g}^{(t+1)} \end{aligned}$$

b)

b)

$$\begin{aligned} \bar{w}_{ix} &= \sum_i \bar{i}^{(t)} \frac{\partial i^{(t)}}{\partial w_{ix}} \\ &= \sum_i \bar{i}^{(t)} \frac{\partial i^{(t)}}{\partial w_{ix}^{(t+1)} + w_{ih}^{(t+1)}} \cdot \frac{\partial w_{ix}^{(t+1)} + w_{ih}^{(t+1)}}{\partial w_{ix}} \\ &= \sum_i \bar{i}^{(t)} \sigma'(w_{ix}^{(t+1)} + w_{ih}^{(t+1)}) \cdot x^{(t)} \end{aligned}$$

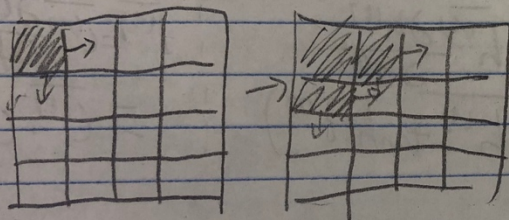
2. a)

2 a) Since we know the hidden unit dimension is H , then according to the function we know that $\dim(W_{in}^{T(i,j)}) = \dim(W_{in}^{T(i-1,j)}) = \dim(W_{in}^{T(i,j-1)}) = H$. According to question we know input dimension is D , then $\dim(x^{(i,j)}) = D$. To let $\dim(W_{in}^{T(i,j)}) = H$, we should let $\dim(W_{in}) = H \times D$. Similarly, we know $\dim(h^{(i,j)}) = \dim(h^{(i,j-1)}) = H$, then $\dim(W_{in}^T) = \dim(W_{in}) = H \times D$. So the total number of weight is $\dim(W_{in}^T) + \dim(W_{in}) + \dim(W_{in}^T) = [H \times D + 2(H \times H)]$

Since we know the total number of weight is $H \times D + 2(H \times H)$, and for the function of $h^{(i,j)}$, each h need to compute $(H \times D + 2(H \times H))$ times, then for each hidden unit it takes $O(H \times D + 2(H \times H))$ time. According to the question, we know that $\dim(\text{grid}) = G \times G$, then for the total time, it should be $O(H \times D + 2(H \times H)) \cdot G \cdot G$
 $= [O((H \times D + 2H^2) G^2)]$

b)

b) It should be $O(2G-1)$ steps. If we know $h^{(i,j)}$, then we can compute its adjacent. If computing $h^{(i,j)}$ takes 1 time, then for all hidden activations, it should be $O(2G-1)$ steps.



c)

Advantage: MDRNNs are more robust to input warping than convolution networks (Multi-Dimensional Recurrent Neural Networks)¹ and capable of modeling long-term sequential dependencies. (Recent Advances in Recurrent Neural Networks)²

Disadvantage: It has higher computational complexity since it uses LSTM. (Recent Advances in Recurrent Neural Networks)

3. a)

3, a)

According to question, we know

$$\begin{aligned}
 S^{(k+1)} &= f(S^{(k)}) \\
 &= (\theta^{(k+1)}, p^{(k+1)}) \\
 &= (\theta^{(k)} + p^{(k)}, p^{(k)} - \alpha \nabla J(\theta^{(k)}))
 \end{aligned}$$

We need $S^{(k)}$, and $S^{(k)} = f^{-1}(S^{(k+1)})$

$$\begin{aligned}
 &= f^{-1}((\theta^{(k+1)}, p^{(k+1)})) \\
 &= (\theta^{(k)}, p^{(k)})
 \end{aligned}$$

According to question we also know $S^{(k)} = (\theta^{(k)}, p^{(k)})$

So, from the function of $\theta^{(k+1)}, p^{(k+1)}$, we can get

$$\begin{aligned}
 \theta^{(k+1)} &= \theta^{(k)} + p^{(k)} \\
 p^{(k)} &= \frac{p^{(k+1)} + \alpha \nabla J(\theta^{(k)})}{\rho}
 \end{aligned}$$

1. <https://arxiv.org/pdf/0705.2011.pdf>
2. <https://arxiv.org/pdf/1801.01078.pdf>

b)

$$b) \det \frac{\partial s^{(k+1)}}{\partial s^{(k)}}$$

According to question, we can get .

$$\frac{\partial p^{(k+1)}}{\partial p^{(k)}} = \beta \quad \frac{\partial \theta^{(k+1)}}{\partial p^{(k)}} = \beta$$

$$\frac{\partial p^{(k+1)}}{\partial \theta^{(k)}} = -\alpha \frac{\partial \nabla J(\theta^{(k)})}{\partial \theta^{(k)}} \quad \frac{\partial \theta^{(k+1)}}{\partial \theta^{(k)}} = 1 - \alpha \frac{\partial \nabla J(\theta^{(k)})}{\partial \theta^{(k)}}$$

If we assume $\frac{\partial \nabla J(\theta^{(k)})}{\partial \theta^{(k)}} = 0$, then

$$\det \frac{\partial s^{(k+1)}}{\partial s^{(k)}} = \det \begin{pmatrix} \frac{\partial p^{(k+1)}}{\partial p^{(k)}} & \frac{\partial p^{(k+1)}}{\partial \theta^{(k)}} \\ \frac{\partial \theta^{(k+1)}}{\partial p^{(k)}} & \frac{\partial \theta^{(k+1)}}{\partial \theta^{(k)}} \end{pmatrix}$$

$$= \det \begin{pmatrix} \beta & 0 \\ \beta & 1 \end{pmatrix}$$

0 is zero matrix

and 1 is identity matrix.

$$= \beta$$