

CSE 891 : Homework 3
Net-id : galibasa

① Binary addition

We can use the 3 hidden units to distinguish between the 3 cases:

i) h_1 to check $(\text{input1} + \text{input2} + \text{carry}) < 1$ or not

ii) h_2 to check $(\text{input1} + \text{input2} + \text{carry}) < 2$ or not
 [if the sum ≥ 2 , we have to carry 1 to the next time step, so this h_2 is crucial for forwarding carry to the next time step]

iii) h_3 to check $(\text{input1} + \text{input2} + \text{carry}) = 3$ or not

$$\text{So, } w_{xh} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}, w_{hh} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix}, b_n = \begin{bmatrix} -0.1 \\ -1.1 \\ -2.1 \end{bmatrix}$$

↳ only consider h_2 to forward carry (see ii)

Then, we have to ~~make~~^{check} sure, when $h_1 = 1$ ($\text{sum} \geq 1$), whether h_2 and h_3 both are same (00 or 11) or not (01 or 10).

$\underbrace{h_2 \ h_3}_{\text{then, } y=1}$

$\underbrace{\downarrow \ \downarrow \ \downarrow \ \downarrow}_{h_2 \ h_3} \quad \underbrace{\downarrow \ \downarrow \ \downarrow \ \downarrow}_{h_2 \ h_3}$
then $y=0$

So,

$$w_{hy} = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \quad b_y = [-0.1]$$

(Ans).

② LSTM Gradient:

$$\begin{aligned}
 ②-1 \quad \overline{h^{(t)}} &= \frac{dL}{di^{(t+1)}} \cdot \frac{di^{(t+1)}}{dh_t^{(t)}} + \frac{dL}{df^{(t+1)}} \cdot \frac{df^{(t+1)}}{dh_t^{(t)}} + \frac{dL}{do^{(t+1)}} \cdot \frac{do^{(t+1)}}{dh_t^{(t)}} \\
 &\approx \underbrace{\dots}_{\text{...}} + \frac{dL}{dg^{(t+1)}} \cdot \frac{dg^{(t+1)}}{dh_t^{(t)}} \\
 &= \overline{i^{(t+1)}} \cdot \underbrace{i^{(t+1)}(1 - i^{(t+1)})}_{\sigma(x) = \sigma(x)(1 - \sigma(x))} \cdot \omega_{hi} + \overline{f^{(t+1)}} \cdot \underbrace{f^{(t+1)}(1 - f^{(t+1)})}_{\sigma'(x) = \dots} \cdot \omega_{hf} \\
 &\quad + \overline{o^{(t+1)}} \cdot \underbrace{o^{(t+1)}(1 - o^{(t+1)})}_{\sigma'(x) = \dots} \cdot \omega_{ho} \\
 &\quad + \overline{g^{(t+1)}} \cdot \underbrace{(1 - \tanh^2(\omega_{xg} x^{(t+1)} + \omega_{hg} h^{(t)}))}_{\tanh(x) = \frac{1 - \tanh^2(x)}{1 + \tanh^2(x)}} \cdot \omega_{hg} \\
 &\quad [+ \frac{dL}{dC^{(t)}} \text{ if loss function has } h^{(t)} \text{ term}]
 \end{aligned}$$

$$\begin{aligned}
 \overline{c^{(t)}} &= \frac{dL}{dh^{(t)}} \cdot \frac{dh^{(t)}}{dc^{(t)}} + \frac{dL}{dc^{(t+1)}} \cdot \frac{dc^{(t+1)}}{dc^{(t)}} \\
 &= \overline{h^{(t)}} \cdot o^{(t)}(1 - \tanh^2(c^{(t)})) + \overline{c^{(t+1)}} \cdot f^{(t+1)}
 \end{aligned}$$

$$\overline{g^{(t)}} = \overline{c^{(t)}} \cdot i^{(t)}$$

$$\overline{o^{(t)}} = \overline{h^{(t)}} \cdot \tanh(c^{(t)})$$

$$\overline{f(t)} = \overline{c^{(t)}} \cdot c^{(t-1)}$$

$$\overline{o^{(t)}} = \overline{c^{(t)}} \cdot g^{(t)}$$

(2-2)

$$\begin{aligned}\overline{\omega_{x_i}} &= \sum_{t=1}^T \frac{\partial L}{\partial i^{(t)}} \cdot \frac{d o^{(t)}}{d \omega_{x_i}} \\ &= \sum_{t=1}^T \overline{c^{(t)}} \cdot \underbrace{o^{(t)} \cdot (1 - o^{(t)})}_{o'(x)} \cdot \overline{x^{(t)}} \\ &\quad \text{or } o'(x) = \sigma(x) \cdot (1 - \sigma(x))\end{aligned}$$

(2-3)

If, $f(t)$ is close to 1 and $\overline{o^{(t)}}$ & $\overline{o^{(t)}}$ are close to 0

then

$$\begin{aligned}\overline{c^{(t)}} &= \overline{h^{(t)}} o^{(t)} (1 - \tanh^2(\overline{c^{(t)}})) + \overline{c^{(t+1)}} \cdot f^{(t+1)} \\ \Rightarrow \overline{c^{(t)}} &\approx \overline{c^{(t+1)}} \quad [\text{as } o^{(t)} \rightarrow 0 \text{ and } f^{(t+1)} \rightarrow 1]\end{aligned}$$

So, gradient does not change in this case.

$$\begin{aligned}\text{Also, } \overline{h^{(t)}} &= \overline{o^{(t+1)}} \cdot o + f^{(t+1)} \cdot o + \overline{o^{(t+1)}} \cdot o \\ &\quad + \overline{g^{(t+1)}} \cdot (1 - \tanh^2 \dots)\end{aligned}$$

if loss has $h^{(t)}$ \Rightarrow $\overline{[+ \frac{\partial L}{\partial h^{(t)}}]}$ [as $\overline{o^{(t+1)}} \approx 0$, $\overline{o^{(t+1)}} \approx 0$, $(1 - f^{(t+1)}) \approx 0$ and $\overline{g^{(t+1)}} = \overline{c^{(t+1)}} \cdot \overline{o^{(t+1)}} = 0$]

So, the gradient does not vanish or explode in that case.

③ ①

$$\text{Conv3-64} : 64 \times 3 \times 3 \times 3 + 64 = 1728 + 64 = \boxed{1792}$$

$$\text{Max-pool} : 0 \quad (224 \rightarrow 112)$$

$$\text{Conv3-128} : 128 \times 64 \times 3 \times 3 + 128 = \boxed{73856}$$

$$\text{Max-pool} : 0 \quad (112 \rightarrow 56)$$

$$\text{Conv3-256} : 256 \times 128 \times 3 \times 3 + 256 = \boxed{295168}$$

$$\text{Conv3-256} : 256 \times 256 \times 3 \times 3 + 256 = \boxed{590080}$$

$$\text{Max-pool} : 0 \quad (56 \rightarrow 28)$$

$$\text{FC-1024} : 1024 \times \cancel{28 \times 28} + 1024 =$$

$$1024 \times 28 \times 28 \times 256 + 1024 = 205520896 + 1024 \\ = \boxed{205521920}$$

$$\text{FC-1000} : 1000 \times 1024 + 1000 = 1025000$$

$$\text{Total} : 207507816$$

③ - ②

The left one is easier to learn in terms of exploding / vanishing gradient.

Because, the left one employs an identity mapping
 $(y = f(x) + x)$

which is very easy to learn and can make the model robust from vanishing / exploding gradient problem.

Also, post-activation in the right can not guarantee that the output will not cause vanishing / exploding gradient.

④ Auto-Regressive Model

④-1 Pixel CNN

①-a) The $O(\cdot)$ of PixelCNN would be same as typical CNN. Because masking in PixelCNN can affect complexity only in constant factors. So, the number of connections

$$\Rightarrow O(dHWK^2)$$

①-b) As we can evaluate each pixel parallelly (in training)

the number of sequential operation would be

$$\rightarrow O(d).$$

But in case of testing, we can not get/evaluate all pixels before getting any all the prior pixels.

So, in that case parallel computation would

not help. and the # of sequential operation

would be :

$$\rightarrow O(HWdK^2)$$

(4-2)

MDRNN

(4-2-a)

The total number of connections:

$$O(d + Wk^2) \quad [\text{as } \# \text{ of input channels} = k \\ \# \text{ of recurrent neurons} = k \\ \# \text{ of layers} = d]$$

(4-2-b)

The number of sequential operations:

$$O(\sqrt{Hwd}) \quad [\text{as here we can do it} \\ \text{partially sequential order,} \\ \text{like if the hidden nodes} \\ \text{from north and west are} \\ \text{evaluated already, we can} \\ \text{evaluate the current node})$$

④ ③

PixelCNN can be parallelized while training. And it captures spatial representation more strongly. Also it ~~uses less~~ focuses less on the sequential pattern which may fail to capture long-term dependencies.

(it has $O(d)$ seq. operations)

MDRNN can not be parallelize like PixelCNN.
it has $O(\sqrt{Hw})$ sequential operations.

But, MDRNN can capture sequential representation more strongly than PixelCNN. Yet, it may lack in learning spatial representation with respect to Pixel CNN.