

①

$$a) V_2 = \begin{bmatrix} -0.185 & 0.245 & -0.059 \\ -0.196 & -0.032 & 0.209 \end{bmatrix}$$

$$S_2 = \begin{bmatrix} 0.0106 & 0.0699 & 0.1099 \\ & & \end{bmatrix}$$

$$S_2 = \begin{bmatrix} 0.0106 & 0.070 & 0.110 \\ 0.141 & 0.170 & 0.021 \end{bmatrix}$$

$$W_2 = \begin{bmatrix} -0.182 & 0.308 & 0.674 \\ 0.502 & -0.858 & -0.934 \end{bmatrix}$$

False

③ a) $h_t = o_t \odot \tanh(c_t)$

$$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o)$$

$x_t = 0$

$$= \sigma[U_o h_{t-1} + b_o]$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \quad \text{--- (3)}$$

$$f_t = \sigma_x[U_f h_{t-1} + b_f] \quad x_t = 0 \quad \text{--- (1)}$$

$$i_t = \sigma[U_i h_{t-1} + b_i] \quad , \quad x_t = 0 \quad \text{--- (2)}$$

$$\tilde{c}_t = \tanh(U_c h_{t-1} + b_c) \quad \text{--- (4)}$$

Putting (1), (2), (4) in (3)

$$c_t = \sigma[U_o h_{t-1} + b_o] \odot c_{t-1} + \sigma[U_i h_{t-1} + b_i] \odot \tanh[U_c h_{t-1} + b_c]$$

$$h_t = \sigma[U_o h_{t-1} + b_o] \odot \left\{ \sigma[U_o h_{t-1} + b_o] \odot c_{t-1} + \sigma[U_i h_{t-1} + b_i] \odot \tanh[U_c h_{t-1} + b_c] \right\}$$

so $h_t \neq h_{t-1}$,

Qualitatively due to input zero also, network may decide to forget something from long term memory as even input to the long term memory

True

3) b) By equation, the propagation of gradient-

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial h_t}{\partial o_t} \times \frac{\partial o_t}{\partial h_{t-1}} + \frac{\partial h_t}{\partial c_t} \times \frac{\partial c_t}{\partial f_t} \times \frac{\partial f_t}{\partial h_{t-1}} + \dots$$

$$\frac{\partial c_t}{\partial c_{t-1}} = f_t$$

~~Star~~ In the long term part of the LSTM, cell $\frac{\partial c_t}{\partial c_{t-1}}$ becomes smaller when f_t is small.

By the concept of forget gate, if the previous input is to be forgotten $f(t)$ should be small, so error does not affect previous states.

True

c) $i_t, f_t = \text{output of sigmoid} [x \text{ sigmoid} \leq 1]$ so non-negative.

for o_t , output gate also sigmoid ~~non~~ output.

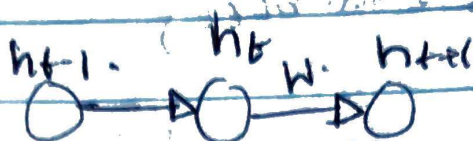
Sigmoid bounded between 0,1 so non-negative output

d) **False**

f_t, i_t, o_t are independent of each other.

f_t, i_t, o_t can be seen as individual probability of forgetting data, adding data and output data, but they do not sum to 1.

a)



$$\frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_{t+1}} \times \frac{\partial h_{t+1}}{\partial h_t}$$

$$h_{t+1} = \text{act}(\cancel{z_t})$$

$$h_{t+1} = \text{act}(\underbrace{z_t}_W)$$

z_t as per question.

$$\frac{\partial h_{t+1}}{\partial h_t} = \sigma'(z_t) W^T$$

$$\boxed{\nabla h_t = \nabla h_{t+1} \times \sigma'(z_t) \times W^T}$$

Dimensional check

$$h_t = 64$$

$$W = 64 \times 64$$

$$(64 \times 32)$$

$$(64 \times 1) \times (1 \times 32)$$

$$(64 \times 32)$$

$$(32 \times 32) \times 64$$

b) $\sigma'(0) = \frac{1}{4}$

$$\nabla h_t = \nabla h_{t+1} \times \frac{1}{4} \times W$$

$$\text{if } \boxed{\frac{W}{\nabla h_t} > 4}$$

$\nabla h_t > 1$, so gradients will explode.

$$\text{if } \boxed{W < 4}$$

$\nabla h_t < 1$, so the gradients will vanish.

$$\boxed{d > 4}$$

Programming Exercises

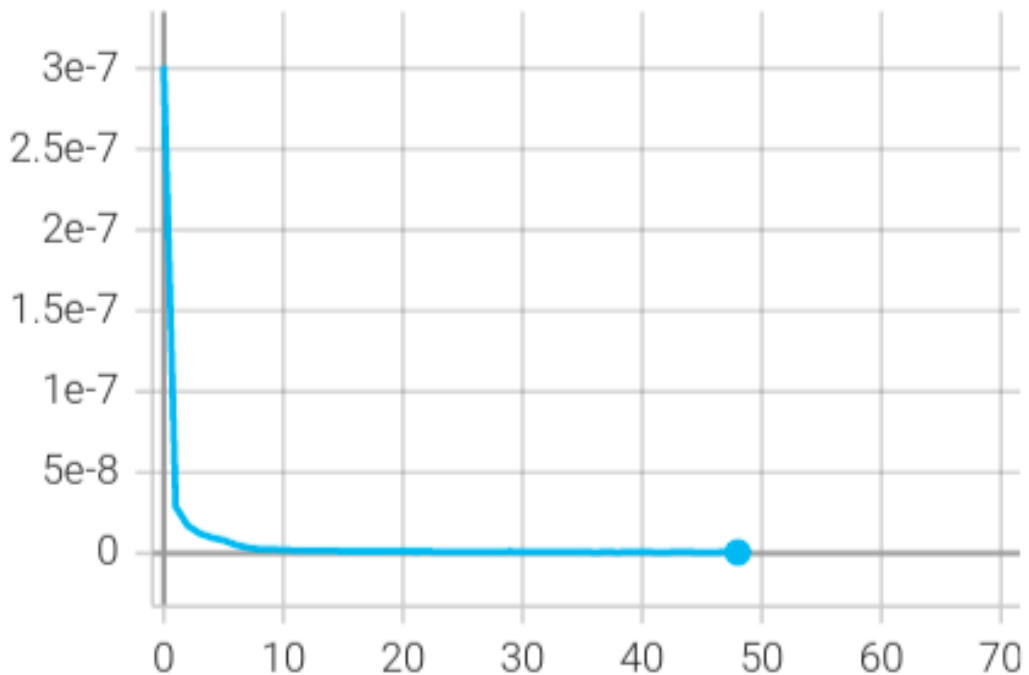
Problem 1

Network structure:

- Three lstm layers, of the following features
 - 1st lstm, Input size 17, hidden size 128
 - 2nd lstm, Input size 128, hidden size 256
 - 3rd lstm, Input size 256, hidden size 512
- One MLP of the following feature
 - Input size 512, output size 17
- Number of trainable parameters: 83985
- Number of iterations: 50
- Learning rate = 0.001
- Optimizer: Adam
- Train loss criterion: MSELoss
- Train loss vs batch number
- Batch size: 256
- Scheduler: Step LR with 0.1 gamma, and 25 epochs step size
 - Train MSELoss vs epoch

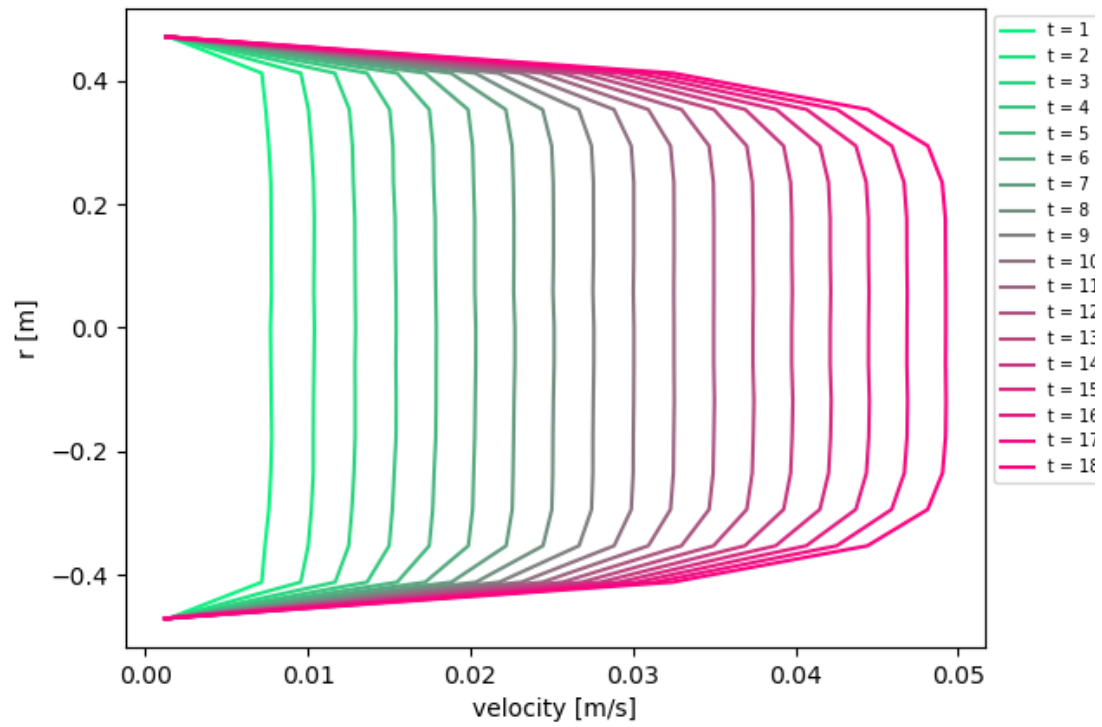
MSE train loss

tag: MSE train loss

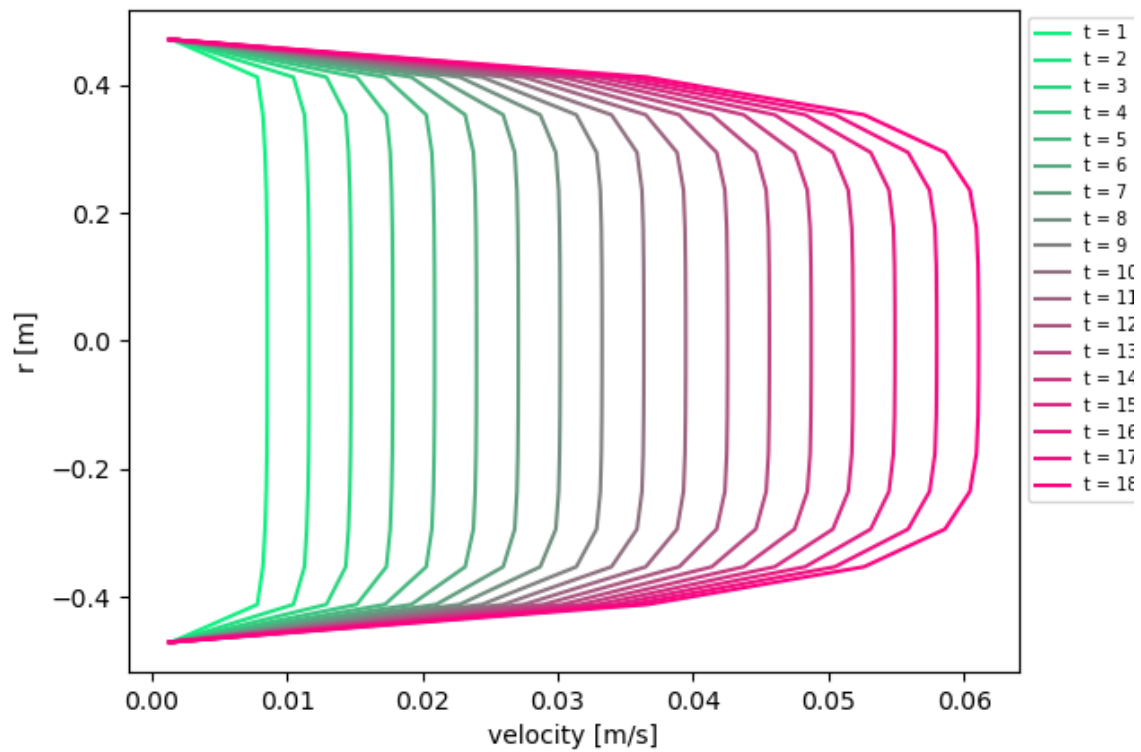


- Final test MSELoss: **0.00028252945048734546**
- Final test L1Loss: **0.06915978819597512**

Results diagram



Prediction



Ground truth