

1.1 Convolutional Neural Networks

(a)

Given by $\lceil (W - F + 2P)/S \rceil + 1$; where W = original dimension; F = filter dimension; S = stride and P = padding

Thus, $\lceil 17/4 \rceil + 1 \times \lceil 7/4 \rceil + 1 = 5 \times 2$

(b)

Given by $F \times \lceil (H-D*(K-1)+2P-1)/S \rceil + 1 \times \lceil (W-D*(K-1)+2P-1)/S \rceil + 1$

(c) (i)

Given by $f_w(x)[r] = \sum_{k=1 \text{ to } k=5} \sum_{i=1 \text{ to } i=3} x[2(r-1)+i][k] W_{1,k,i}$

Dimension is R^3

(ii)

Dimension of $\partial f_w(x)/\partial W$ is $R^{(3) \times (1 \times 5 \times 3)}$

$(\partial f_w(x)/\partial W)[r, c, k, i] = x[2(r-1)+i][k]$

(iii)

Dimension of $\partial f_w(x)/\partial x$ is $R^{(3) \times (5 \times 7)}$

$(\partial f_w(x)/\partial x)[r, k, i] = W_{1,k,i-2(r-1)}$, when $i-2(r-1) \in [1, 3]$
= 0, otherwise

(iv)

Dimension is $R^{1 \times 5 \times 3}$

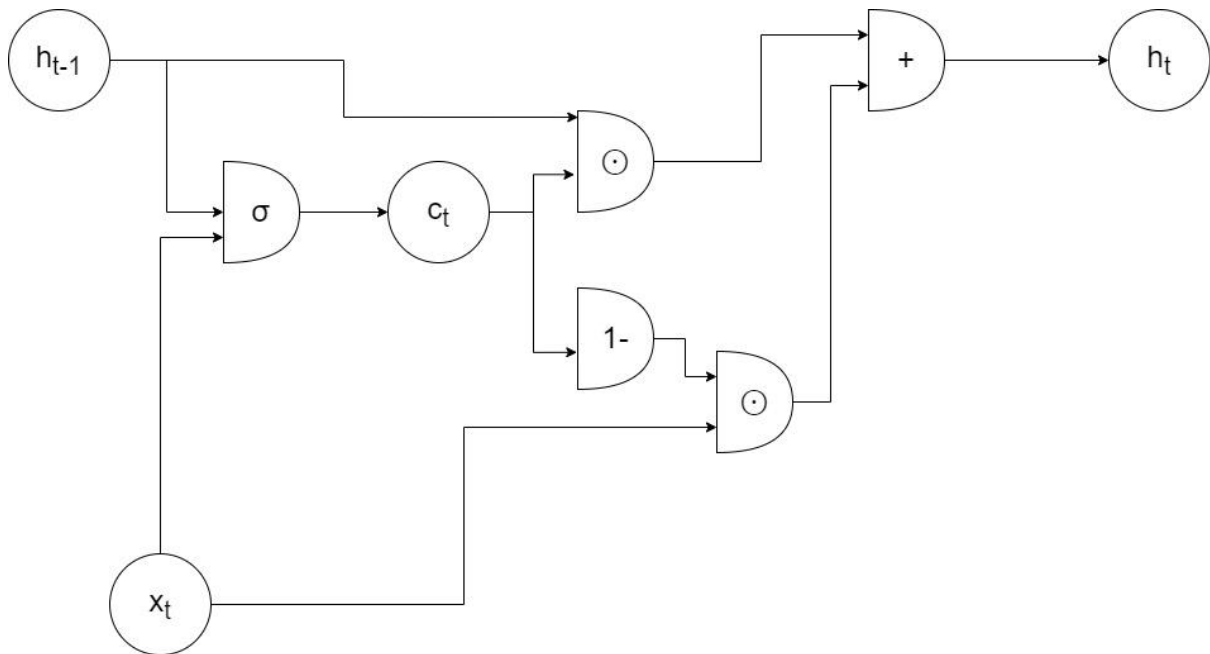
$(\partial l/\partial W)[1, k, i] = \sum_{r=1 \text{ to } r=2} (\partial l/\partial f_w(x))[r] x[2(r-1)+i, k]$

The similarities are that both forward and backward pass involve a convolution. The differences are in the backward pass, because stride becomes dilation.

1.2 Recurrent Neural Networks

1.2.1 Part 1

(a)



(b)

The dimension of $c[t]$ is the same as the dimension of $W_c x[t]$, which is R^m

(c)

$$\frac{\partial l}{\partial W_x} = \left(\frac{\partial l}{\partial h[t]} \right) \cdot \left(\frac{\partial h[t]}{\partial W_x} \right) \quad \{\text{Chain Rule}\}$$

$$\text{Now, } \frac{\partial h[t]}{\partial W_x} = (1 - c[t]) \cdot x[t]$$

$$\text{Thus, } \frac{\partial l}{\partial W_x} = \left(\frac{\partial l}{\partial h[t]} \right) \cdot (1 - c[t]) \cdot x[t].$$

$$\frac{\partial l}{\partial W_x} \text{ has the same dimension as } W_x, \text{ which is } R^{(m \times n)}$$

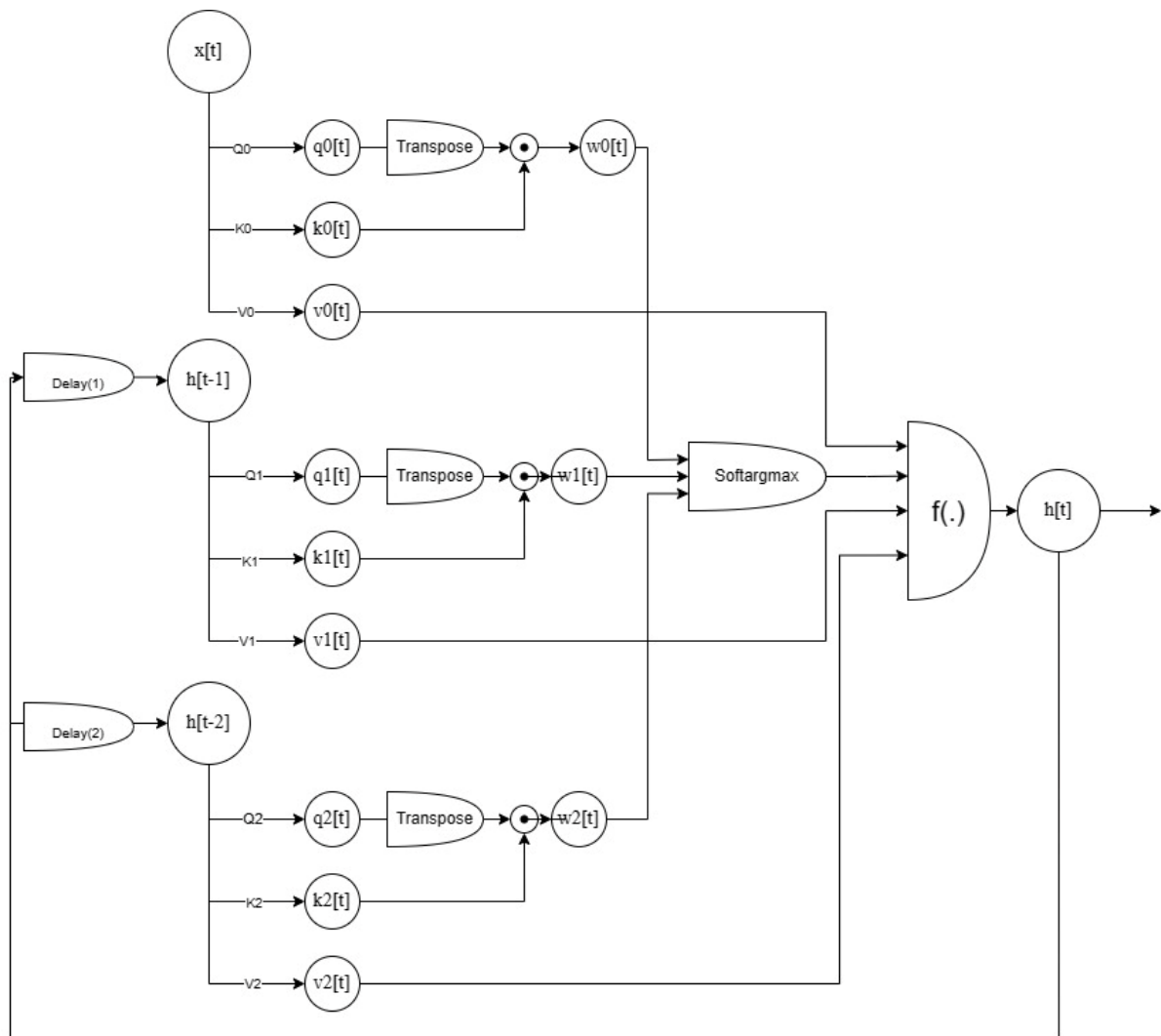
The similarity between the forward and backward pass is that the same computational graph is used but with the inputs and outputs swapped. In the forward pass, we use $x[1], \dots, x[t]$ to calculate $h[1], h[2], \dots, h[t]$ recurrently and in the backward pass we use $\frac{\partial l}{\partial h[t]}$ to calculate the input gradients $\frac{\partial l}{\partial x[t]}$ recurrently.

(d)

This RNN cannot be subject to exploding gradients because the state vector h_t is not multiplied by any matrix from one timestep to the next. It can be subject to vanishing gradients because on each iteration the state vector h_t is multiplied with c_t , which is a vector of values between 0 and 1.

1.2.2 Part 2

(a)



(b)

\mathbb{R}^3

(c)

$$q_0[t], q_1[t], q_2[t] \dots q_k[t] = Q_0 x[t], Q_1 h[t-1], Q_2 h[t-2] \dots Q_k h[t-k]$$

$$k_0[t], k_1[t], k_2[t] \dots k_k[t] = K_0 x[t], K_1 h[t-1], K_2 h[t-2] \dots K_k h[t-k]$$

$$v_0[t], v_1[t], v_2[t] \dots v_k[t] = V_0 x[t], V_1 h[t-1], V_2 h[t-2] \dots V_k h[t-k]$$

$$w_i[t] = q_i[t]^T k_i[t]$$

$$a[t] = \text{softargmax}([w_0[t], w_1[t], w_2[t] \dots w_k[t]])$$

$$h[t] = \sum (i = 0 \text{ to } i = k) a_i[t] v_i[t]$$

(d)

We can share weights for all the Q, K and V

Thus,

$$q_0[t], q_1[t], q_2[t] \dots q_{t-1}[t] = Q_0 x[t], Q_h[t-1], Q_h[t-2] \dots Q_h[1]$$

$$k_0[t], k_1[t], k_2[t] \dots k_{t-1}[t] = K_0 x[t], K_h[t-1], K_h[t-2] \dots K_h[1]$$

$$v_0[t], v_1[t], v_2[t] \dots v_{t-1}[t] = V_0 x[t], V_h[t-1], V_h[t-2] \dots V_h[1]$$

$$w_i[t] = q_i[t]^T k_i[t]$$

$$a[t] = \text{softmax}([w_0[t], w_1[t], w_2[t] \dots w_{t-1}[t]])$$

$$h[t] = \sum (i = 0 \text{ to } i = t-1) a_i[t] v_i[t]$$

(e)

$$\partial h[t] / \partial h[t-1] = \partial (a_0[t] v_0[t] + a_1[t] v_1[t] + a_2[t] v_2[t]) / \partial h[t-1] \quad [1]$$

Where,

$$a_0[t] = \exp(q_0[t]^T k_0[t]) / \sum \exp(q_i[t]^T k_i[t])$$

$$a_1[t] = \exp(q_1[t]^T k_1[t]) / \sum \exp(q_i[t]^T k_i[t])$$

$$a_2[t] = \exp(q_2[t]^T k_2[t]) / \sum \exp(q_i[t]^T k_i[t])$$

Now,

$$\partial (a_0[t] v_0[t]) / \partial h[t-1] = -v_0[t] (k_1 Q_1 + q_1^T K_1) (\exp(q_0[t]^T k_0[t])) (\exp(q_1[t]^T k_1[t])) / (\sum \exp(q_i[t]^T k_i[t]))^2 \quad [2]$$

Similarly,

$$\partial (a_2[t] v_2[t]) / \partial h[t-1] = -v_2[t] (k_1 Q_1 + q_1^T K_1) (\exp(q_2[t]^T k_2[t])) (\exp(q_1[t]^T k_1[t])) / (\sum \exp(q_i[t]^T k_i[t]))^2 \quad [3]$$

Finally,

$$\partial (a_1[t] v_1[t]) / \partial h[t-1] = v_1[t] \partial a_1[t] / \partial h[t-1] + a_1[t] \partial v_1[t] / \partial h[t-1] \quad [4]$$

Now,

$$\partial v_1[t] / \partial h[t-1] = V_1 \quad [5]$$

$$\begin{aligned} \partial a_1[t] / \partial h[t-1] = & \{ (k_1 Q_1 + q_1^T K_1) (\exp(q_1[t]^T k_1[t])) / (\sum \exp(q_i[t]^T k_i[t])) - \\ & \{ (\exp(q_1[t]^T k_1[t])) (k_1 Q_1 + q_1^T K_1) (\exp(q_1[t]^T k_1[t])) \} / (\sum \exp(q_i[t]^T k_i[t]))^2 \end{aligned} \quad [6]$$

We substitute [5] and [6] in [4]. And then [2], [3] and [4] in [1] to get the final expression.

(Not written here in one go for the sake of clarity)

(f)

$$\sum_{i=1}^k (\partial h[i+T] / \partial h[T]) (\partial l / \partial h[i+T])$$

1.3 Debugging loss curves

(1)

As discussed in the lecture, this may be happening due to exploding gradients.

(2)

This is possible when the labels are being incorrectly assigned. Initially when the weights are randomly initialized the loss is given by $\ln(\text{number of classes})$, however if say the labels are swapped completely (i.e in a binary classifier 1s are predicted as 0s and 0s are predicted as 1s) then the model loss and accuracy will be worse than a randomly initialized model. Such a situation may arise due to several reasons related to the training process and hyperparameters.

(3)

Some ways to fix them is using regularization techniques and gradient clipping.

(4)

$$\text{Initial loss} = \ln(4) = 1.386$$

$$\text{Initial accuracy} = \frac{1}{4} = 0.25$$