

DL Lab4: Back-Propagation Through Time

0756079 陳冠聞

May 8, 2019

1 Introduction

Recurrent neural networks is a specific class of artificial neural network. In traditional feedforward network, input examples are fed to the network and transformed into an output. On the other hand, recurrent neural networks take as their input not just the current input example they see, but also what they have perceived previously in time. This allows it to model the temporal relationship in input data, and it has many application such as speech recognition, machine translation and so on.

In this lab, I will implement the recurrent neural network in python using only *NumPy* for matrix operation, others function are all implemented by myself.

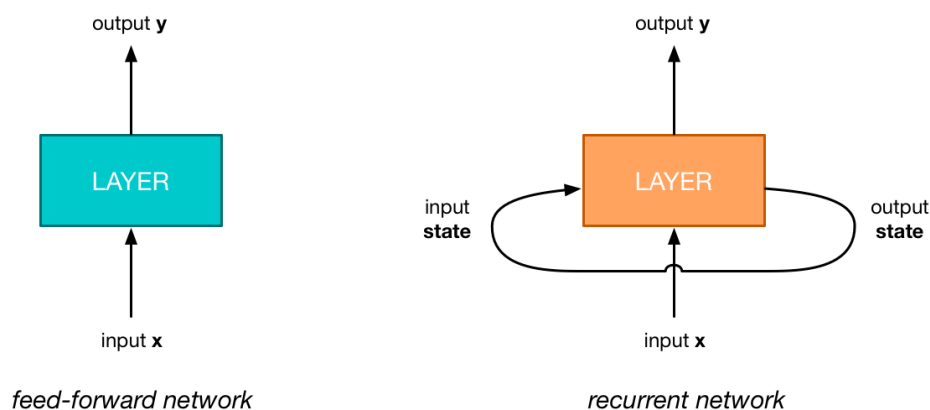


Figure 1: Feedforward v.s. Recurrent Neural Network

2 Recurrent Neural Network

Recurrent neural network is a sequence to sequence mapping between input sequence $x^{(1)}, x^{(2)}, \dots, x^{(\tau)}$ to the target sequence $y^{(1)}, y^{(2)}, \dots, y^{(\tau)}$. In this section, I will explain how the network generate output through forward propagation, and how it can learn from data through back-propagation.

2.1 Forward Propagation

Let the weights of input-to-input, hidden-to-hidden and hidden-to-output connection be U , V and W respectively. In addition, let the hidden state and output unit at timestep t be the $h^{(t)}$ and $o^{(t)}$, the forward propagation at timestep t can be formulated as following

$$\begin{aligned} a^{(t)} &= b + Wh^{(t-1)} + Ux^{(t)} \\ h^{(t)} &= \tanh(a^{(t)}) \\ o^{(t)} &= c + Vh^{(t)} \end{aligned} \tag{1}$$

where b and c are the bias term of hidden unit and output unit, and $a^{(t)}$ is value of hidden unit before activation function \tanh . The diagram of the forward propagation can refer to Figure 2, and the implementation code is shown in Figure 5.

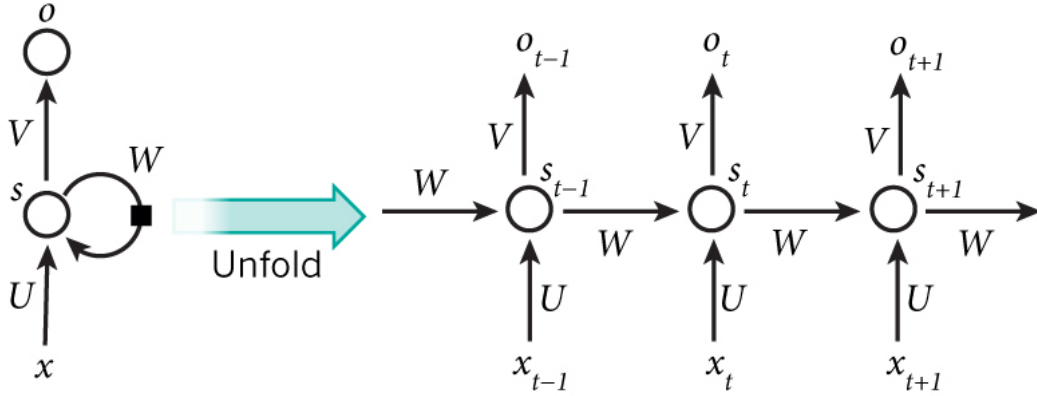


Figure 2: Forward Propagation of RNN

2.2 Back Propagation Through Time

Given the corresponding input sequences and target sequences, we want to learn the model parameters $\theta = \{W, U, V, b, c\}$ such that the output of the learned network given the input sequence will close to the target sequence. We can use the same concept of back-propagation in the feedforward network, but the difference is that the recurrent network has many timesteps, so we have to back-propagate through time like the unfold graph in Figure 2.

The gradient of each model parameter is shown below, and the detail mathematical derivation is left in appendix for the conciseness.

$$\begin{aligned}\nabla_c L &= \sum_t \left(\frac{\partial L}{\partial o^t} \right) \\ \nabla_b L &= \sum_t H^{(t)} \left(\frac{\partial L}{\partial h^t} \right) \\ \nabla_V L &= \sum_t \left(\frac{\partial L}{\partial o^t} \right) (h^{(t)})^T \\ \nabla_U L &= \sum_t H^{(t)} \left(\frac{\partial L}{\partial h^t} \right) (x^{(t)})^T \\ \nabla_W L &= \sum_t H^{(t)} \left(\frac{\partial L}{\partial h^t} \right) (h^{(t-1)})^T\end{aligned}\tag{2}$$

3 Implementation

3.1 Data Generation

For data generation, I use the *numpy.random* to generate two 8-bits numbers x_1 and x_2 , then I follow the the binary addition procedure to get the correct answer y . When fed to the network, the input at timestep t is consist of two corresponding bit in x_1 and x_2 , and from the least significant bit to the most significant bit.

```
class BinaryAdditionDataGenerator():
    def generate_data(self, n_bits=8):
        x1 = np.random.randint(2, size=n_bits)
        x2 = np.random.randint(2, size=n_bits)
        y = [0 for i in range(n_bits)]
        carry = 0
        for i, (a, b) in enumerate(zip(x1, x2)):
            y[i] = a + b + carry
            carry = 0
            if y[i] >= 2:
                y[i] %= 2
                carry = 1
        x = np.array([x1, x2]).T
        y = np.array(y)
        return x, y
```

Figure 3: Code Snippet for Data Generation

3.2 Nerual Net

First, for the neural net architecture, I initialize the weights of W, U, V using **Xavier initialization** and set the bias term b and c to all zeros.

```
class RNN():
    def __init__(self, in_dims=2, hidden_dims=16, out_dims=1):
        # Model dimensions
        self.in_dims = in_dims
        self.hidden_dims = hidden_dims
        self.out_dims = out_dims

        # Model Weights
        self.U = np.random.normal(size=(hidden_dims, in_dims)) * np.sqrt(2/(in_dims+hidden_dims))
        self.V = np.random.normal(size=(out_dims, hidden_dims)) * np.sqrt(2/(hidden_dims+out_dims))
        self.W = np.random.normal(size=(hidden_dims, hidden_dims)) * np.sqrt(2/(hidden_dims+hidden_dims))
        self.b = np.zeros((hidden_dims, 1))
        self.c = np.zeros((out_dims, 1))
        self.zero_grad()

    def zero_grad(self):
        # Set Gradients to zeros
        self.gradient_c = np.zeros((self.c.shape))
        self.gradient_b = np.zeros((self.b.shape))
        self.gradient_V = np.zeros((self.V.shape))
        self.gradient_U = np.zeros((self.U.shape))
        self.gradient_W = np.zeros((self.W.shape))
        self.batch_size = 0

    def update_weight(self, lr=0.1):
        # Update Model parameters
        self.c -= lr * (1/self.batch_size) * self.gradient_c
        self.b -= lr * (1/self.batch_size) * self.gradient_b
        self.V -= lr * (1/self.batch_size) * self.gradient_V
        self.U -= lr * (1/self.batch_size) * self.gradient_U
        self.W -= lr * (1/self.batch_size) * self.gradient_W
        self.zero_grad()
```

Figure 4: Code Snippet for RNN (Init)

Second, for the forward propagation, I follow the Equation 1. I assume the timestep of input sequence start from $t = 1$, and set the dummy variables for timestep $t = 0$ to be consistent with the initial hidden state $h^{(0)}$. In the forward propagation procedure, I record all the intermediate results to do back-propagation.

```
def forward(self, xs):
    T = len(xs)
    # Record from t=0 to t=T (t=0 is initial state)
    self.x = [None for t in range(0, T+1)]
    self.a = [None for t in range(0, T+1)]
    self.h = [None for t in range(0, T+1)]
    self.o = [None for t in range(0, T+1)]
    self.y = [None for t in range(0, T+1)]

    # Foward Pass
    self.h[0] = np.zeros((self.hidden_dims, 1))
    for i, x in enumerate(xs):
        t = i+1
        self.x[t] = x.reshape(-1, 1)
        self.a[t] = self.b + np.matmul(self.W, self.h[t-1]) + np.matmul(self.U, self.x[t])
        self.h[t] = tanh(self.a[t])
        self.o[t] = self.c + np.matmul(self.V, self.h[t])
        self.y[t] = sigmoid(self.o[t])
    return self.y
```

Figure 5: Code Snippet for RNN (Forward Pass)

Lastly, for the back-propagation, I follow the Equation 2. I compute the gradient of model parameters, and accumulate the gradient, then update the model weights after a mini-batch to achieve batch-training.

```
def backward(self, y_true, y_pred):
    T = len(self.h) - 1

    # Compute dL/do
    d_o = [None for t in range(0, T+1)]
    for t in range(1, T+1):
        i = t - 1
        d_o[t] = y_pred[t] - y_true[i]

    # Compute H^*(t)
    H = [None for t in range(0, T+1)]
    for t in range(1, T+1):
        H[t] = np.zeros((self.hidden_dims, self.hidden_dims))
        for d in range(H[t].shape[0]):
            H[t][d, d] = 1-(self.h[t][d]**2)

    # Compute dL/dh
    d_h = [None for t in range(0, T+1)]
    d_h[-1] = np.matmul(self.V.T, d_o[-1])
    for t in range(T-1, 0, -1):
        d_h[t] = np.matmul(self.W.T, np.matmul(H[t+1], d_h[t+1])) + np.matmul(self.V.T, d_o[t])

    # Compute gradient with respect to model parameters {c, b, V, U, W}
    for t in range(1, T+1):
        self.gradient_c += d_o[t]
        self.gradient_b += np.matmul(H[t], d_h[t])
        self.gradient_V += np.matmul(d_o[t], self.h[t].T)
        self.gradient_U += np.matmul(np.matmul(H[t], d_h[t]), self.x[t].T)
        self.gradient_W += np.matmul(np.matmul(H[t], d_h[t]), self.h[t-1].T)
    self.batch_size += 1
```

Figure 6: Code Snippet for RNN (Back-Propagation)

4 Result

In this lab, I train the network using stochastic gradient descent with batch size 8 as well as learning rate 0.005. I use two criterions to evaluate the model performance. One is average **bit error**, which is the number of different bits between the model prediction and the label. Another is the **accuracy**, which considers the prediction is correct if it is exactly the same as the label for all bits, then compute the number of correct divided by the number of data.

From the Figure 7, we can see that the bit error drop very quickly, and accuracy improve dramatically. In addition, the model can achieve 100% accuracy and zero bit error after 3000 iteration.

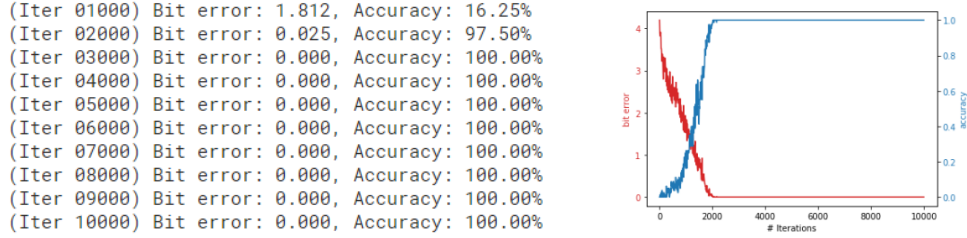


Figure 7: Traing Trend

To test the correctness of model, I generate 2 random integers between $[0, 127]$, and feed them to the network. The result shown in Figure 8 shows that the model successfully learn the task of binary addition.


```
import random
x1 = random.randint(0, 127)
x2 = random.randint(0, 127)
y_pred = net.predict_number(x1, x2)
print("input : (%d, %d)"%(x1, x2))
print("prediction : %d"%(y_pred))
```

```
input : (41, 96)
prediction : 137
```

Figure 8: Result of Adding Two Random Integers

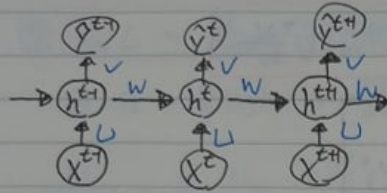
5 Appendix

Belows are the mathematical derivation for the **Back-Propagation Through Time** (BPTT), note that the symbol $H^{(t)}$ is replaced by $D^{(t)}$ in the following figures.

BPTT

No. 075019
Date 陣易明

Notation



$\hat{y}^{(t)}$, predicted label at timestep t

$h^{(t)}$, hidden state at timestep t

$x^{(t)}$, input data at timestep t

$$\begin{cases} a^{(t)} = b + w h^{(t-1)} + U x^{(t)} \\ h^{(t)} = \tanh(a^{(t)}) \end{cases} \quad \begin{cases} d^{(t)} = c + V h^{(t)} \\ \hat{y}^{(t)} = \text{softmax}(d^{(t)}) \end{cases}$$

Loss function $L = L^{(1)} + L^{(2)} + \dots + L^{(T)}$

其中 $L^{(t)} = -y^{(t)} \log \hat{y}^{(t)}$

Derivation

① $\frac{\partial L}{\partial \hat{y}^{(t)}} = 1$

② $(\nabla_{\hat{y}^{(t)}} L)_i = \frac{\partial L}{\partial \hat{y}_i^{(t)}} = \frac{\hat{y}_i^{(t)}}{y_i^{(t)}} - y_i^{(t)}$, $\frac{\partial L}{\partial d^{(t)}} = \hat{y}^{(t)} - y^{(t)}$

(詳見 Softmax + Cross-Entropy 推导)

③ $\nabla_{H^{(t)}} L = \frac{\partial L}{\partial H^{(t)}} = V^T \frac{\partial L}{\partial d^{(t)}}$ (T 为最后一个 time-step)

因 $d_i^{(t)} = \sum_j V_{ij} h_j^{(t)} + c_i$

$\frac{\partial L}{\partial h_j^{(t)}} = \sum_i \frac{\partial L}{\partial d_i^{(t)}} \frac{\partial d_i^{(t)}}{\partial h_j^{(t)}} = \sum_i \frac{\partial L}{\partial d_i^{(t)}} (V_{ij})$

$\therefore \frac{\partial L}{\partial H^{(t)}} = \begin{bmatrix} \sum_i \frac{\partial L}{\partial d_i^{(t)}} V_{i1} \\ \vdots \\ \sum_i \frac{\partial L}{\partial d_i^{(t)}} V_{im} \end{bmatrix} = \begin{bmatrix} V_{11} & \dots & V_{m1} \\ \vdots & \ddots & \vdots \\ V_{1m} & \dots & V_{mm} \end{bmatrix} \begin{bmatrix} \frac{\partial L}{\partial d_1^{(t)}} \\ \vdots \\ \frac{\partial L}{\partial d_m^{(t)}} \end{bmatrix} = V^T \left(\frac{\partial L}{\partial d^{(t)}} \right)$

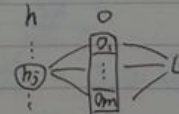


Figure 9: BPTT Derivation (part1)

Derivation

④ $\nabla_{H^{(t)}} L = \frac{\partial L}{\partial H^{(t)}} = W^T \left(\frac{\partial L}{\partial H^{(t+1)}} \right) \nabla^{(t+1)} + V^T \frac{\partial L}{\partial o^{(t)}}$

其中 $\nabla^{(t)} = \begin{bmatrix} (1-h_1^{(t)})^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & (1-h_n^{(t)})^2 \end{bmatrix}$

因 $\begin{cases} a_{\lambda}^{(t+1)} = \sum_j w_{ij} h_j + \square \\ h_{\lambda}^{(t+1)} = \tanh(a_{\lambda}^{(t+1)}) \end{cases}$, 且 $\begin{cases} \tanh(z) = \\ \tanh'(z) = 1 - \tanh^2(z) \end{cases}$

$\frac{\partial L}{\partial h_j^{(t)}} = \sum_{\lambda} \frac{\partial L}{\partial h_j^{(t+1)}} \frac{\partial h_j^{(t+1)}}{\partial a_{\lambda}^{(t+1)}} \frac{\partial a_{\lambda}^{(t+1)}}{\partial h_j^{(t)}} + \sum_{\lambda} \frac{\partial L}{\partial o_{\lambda}^{(t)}} \frac{\partial o_{\lambda}^{(t)}}{\partial h_j^{(t)}}$ 报错了

$= \sum_{\lambda} \left(\frac{\partial L}{\partial h_j^{(t+1)}} \right) (1-h_{\lambda}^{(t+1)})^2 w_{\lambda j} + \sum_{\lambda} \frac{\partial L}{\partial o_{\lambda}^{(t)}} v_{\lambda j}$

$\therefore \frac{\partial L}{\partial H^{(t)}} = \begin{bmatrix} \frac{\partial L}{\partial h_1^{(t)}} \\ \vdots \\ \frac{\partial L}{\partial h_n^{(t)}} \end{bmatrix} = W^T \nabla^{(t+1)} \left(\frac{\partial L}{\partial H^{(t+1)}} \right) + V^T \left(\frac{\partial L}{\partial o^{(t)}} \right)$ ↑ preparing part

⑤ $\nabla_C L = \frac{\partial L}{\partial C} = \sum_{\lambda} \frac{\partial L}{\partial o_{\lambda}^{(t)}}$

因 $o_{\lambda}^{(t)} = C_{\lambda}^{(t)} + \square \Rightarrow \frac{\partial o_{\lambda}^{(t)}}{\partial C_{\lambda}^{(t)}} = 1$

$\frac{\partial L}{\partial C_{\lambda}^{(t)}} = \sum_{\lambda} \frac{\partial L}{\partial o_{\lambda}^{(t)}} \frac{\partial o_{\lambda}^{(t)}}{\partial C_{\lambda}^{(t)}} = \sum_{\lambda} \frac{\partial L}{\partial o_{\lambda}^{(t)}}$

$\therefore \frac{\partial L}{\partial C} = \frac{\partial L}{\partial o^{(t)}}$

⑥ $\nabla_b L = \frac{\partial L}{\partial b} = \sum_{\lambda} \nabla^{(t)} \frac{\partial L}{\partial h_{\lambda}^{(t)}}$

因 $a_{\lambda}^{(t)} = b_{\lambda} + \square \Rightarrow \frac{\partial a_{\lambda}^{(t)}}{\partial b_{\lambda}} = 1$

$\frac{\partial L}{\partial b_{\lambda}} = \sum_{\lambda} \frac{\partial L}{\partial h_{\lambda}^{(t)}} \frac{\partial h_{\lambda}^{(t)}}{\partial a_{\lambda}^{(t)}} \frac{\partial a_{\lambda}^{(t)}}{\partial b_{\lambda}} = \sum_{\lambda} \frac{\partial L}{\partial h_{\lambda}^{(t)}} \cdot (1-h_{\lambda}^{(t)})^2 \cdot 1$

$\therefore \frac{\partial L}{\partial b} = \sum_{\lambda} \nabla^{(t)} \frac{\partial L}{\partial h_{\lambda}^{(t)}}$

Figure 10: BPTT Derivation (part2)

BPTT (2)

Derivation

$$\textcircled{7} \quad \nabla_{\mathbf{V}} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \mathbf{V}} = \sum_t \left(\frac{\partial \mathcal{L}}{\partial \mathbf{a}_t^{(e)}} \right) (\mathbf{h}_t^{(e)})^T$$

$$\textcircled{8} \quad \mathbf{a}_i = \sum_j \mathbf{V}_{ij} \mathbf{h}_j \Rightarrow \frac{\partial \mathbf{a}_i}{\partial \mathbf{V}_{ij}} = \mathbf{h}_j$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{V}_{ij}} = \sum_t \frac{\partial \mathcal{L}}{\partial \mathbf{a}_t^{(e)}} \frac{\partial \mathbf{a}_t^{(e)}}{\partial \mathbf{V}_{ij}} = \sum_t \frac{\partial \mathcal{L}}{\partial \mathbf{a}_t^{(e)}} (\mathbf{h}_j)$$

$$\therefore \frac{\partial \mathcal{L}}{\partial \mathbf{V}} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{V}_{11}} & \dots & \frac{\partial \mathcal{L}}{\partial \mathbf{V}_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mathcal{L}}{\partial \mathbf{V}_{m1}} & \dots & \frac{\partial \mathcal{L}}{\partial \mathbf{V}_{mn}} \end{bmatrix} = \begin{bmatrix} \sum_t \frac{\partial \mathcal{L}}{\partial \mathbf{a}_t^{(e)}} (\mathbf{h}_1) & \dots & \sum_t \frac{\partial \mathcal{L}}{\partial \mathbf{a}_t^{(e)}} (\mathbf{h}_n) \\ \sum_t \frac{\partial \mathcal{L}}{\partial \mathbf{a}_t^{(e)}} (\mathbf{h}_1) & & \sum_t \frac{\partial \mathcal{L}}{\partial \mathbf{a}_t^{(e)}} (\mathbf{h}_n) \end{bmatrix}$$

$$= \sum_t \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{a}_t^{(e)}} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial \mathbf{a}_t^{(e)}} \end{bmatrix} [\mathbf{h}_1^{(e)} \dots \mathbf{h}_n^{(e)}] = \sum_t \left(\frac{\partial \mathcal{L}}{\partial \mathbf{a}_t^{(e)}} \right) (\mathbf{h}_t^{(e)})^T$$

$$\textcircled{8} \quad \nabla_{\mathbf{U}} \mathcal{L} = \sum_t \nabla^{(e)} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t^{(e)}} \right) (\mathbf{x}_t^{(e)})^T$$

$$\textcircled{9} \quad \mathbf{a}_i = \sum_j \mathbf{U}_{ij} \mathbf{x}_j + \square \Rightarrow \frac{\partial \mathbf{a}_i}{\partial \mathbf{U}_{ij}} = \mathbf{x}_j$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{U}_{ij}} = \sum_t \frac{\partial \mathcal{L}}{\partial \mathbf{h}_t^{(e)}} \frac{\partial \mathbf{h}_t^{(e)}}{\partial \mathbf{a}_t^{(e)}} \frac{\partial \mathbf{a}_t^{(e)}}{\partial \mathbf{U}_{ij}} = \sum_t \left(\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t^{(e)}} \right) \cdot (1 - \mathbf{h}_i^{(e)})^2 \cdot (\mathbf{x}_j^{(e)})$$

$$\therefore \frac{\partial \mathcal{L}}{\partial \mathbf{U}} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{U}_{11}} & \dots & \frac{\partial \mathcal{L}}{\partial \mathbf{U}_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mathcal{L}}{\partial \mathbf{U}_{m1}} & \dots & \frac{\partial \mathcal{L}}{\partial \mathbf{U}_{mn}} \end{bmatrix} = \begin{bmatrix} \sum_t \left(\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t^{(e)}} \right) (1 - \mathbf{h}_1^{(e)})^2 (\mathbf{x}_1) & \dots & \\ \vdots & \ddots & \vdots \end{bmatrix}$$

$$= \sum_t \begin{bmatrix} (1 - \mathbf{h}_1^{(e)})^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & (1 - \mathbf{h}_n^{(e)})^2 \end{bmatrix} \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{h}_t^{(e)}} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial \mathbf{h}_t^{(e)}} \end{bmatrix} [\mathbf{x}_1^{(e)} \dots \mathbf{x}_n^{(e)}]$$

$$= \sum_t \nabla^{(e)} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t^{(e)}} \right) (\mathbf{x}_t^{(e)})^T$$

Figure 11: BPTT Derivation (part3)

Derivation

⑨
$$\nabla_W L = \frac{\partial L}{\partial W} = \sum_i \nabla^{(i)} \left(\frac{\partial L}{\partial A^{(i)}} \right) (h^{(i)})^T$$

⊗
$$a_i^{(i)} = \sum_j w_{ij} h_j^{(i-1)} + \square \Rightarrow \frac{\partial a_i^{(i)}}{\partial w_{ij}} = h_j^{(i-1)}$$

$$\frac{\partial L}{\partial w_{ij}} = \sum_i \frac{\partial L}{\partial h_i^{(i)}} \frac{\partial h_i^{(i)}}{\partial a_i^{(i)}} \frac{\partial a_i^{(i)}}{\partial w_{ij}}$$

$$= \sum_i \frac{\partial L}{\partial h_i^{(i)}} (1 - h_i^{(i)})^2 h_j^{(i-1)}$$

$$\therefore \frac{\partial L}{\partial W} = \begin{bmatrix} \frac{\partial L}{\partial w_{11}} & \dots & \frac{\partial L}{\partial w_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial L}{\partial w_{m1}} & \dots & \frac{\partial L}{\partial w_{mn}} \end{bmatrix} = \begin{bmatrix} \sum_i \frac{\partial L}{\partial h_i^{(i)}} (1 - h_i^{(i)})^2 h_i^{(i-1)} & \dots \\ \vdots & \ddots \end{bmatrix}$$

$$= \sum_i \begin{bmatrix} (1 - h_1^{(i)})^2 & 0 \\ \vdots & \vdots \\ 0 & (1 - h_n^{(i)})^2 \end{bmatrix} \begin{bmatrix} \frac{\partial L}{\partial h_1^{(i)}} \\ \vdots \\ \frac{\partial L}{\partial h_n^{(i)}} \end{bmatrix} \begin{bmatrix} h_1^{(i-1)} & \dots & h_n^{(i-1)} \end{bmatrix}$$

$$= \sum_i \nabla^{(i)} \left(\frac{\partial L}{\partial A^{(i)}} \right) (h^{(i-1)})^T$$

Figure 12: BPTT Derivation (part4)