

Chapter 11

RECURRENT NEURAL NETWORK

CONTENTS

11.1	More on Matrix Calculus	521
11.2	Process in RNN	525
11.2.1	Forward Propagation of a Temporal Dataset	525
11.2.2	An Illustrative Example of RNN	527
11.2.3	Backpropagation through Time	529
11.3	Problems with RNN	532
11.4	Gated RNNs	534
11.4.1	Backpropagation through Time	535
11.4.2	Example of Gated RNN: LSTM	537
11.4.3	Uniform Credit Assignment	541

Recurrent neural networks (RNNs) is particularly of use in temporally sequential data alternative to the usual time series data analysis so that RNNs can either classify the entire sequence or generate the subsequent data. The input features are an ordered list of vectors \mathbf{x}_t , for $t = 0, 1, \dots$, and the main difference, in terms of the network structure, from the usual MLP is that, some outputs of the hidden layers will be fed backwardly as the inputs for itself or for previous layers. RNNs are often used in text mining and speech recognition because sentences of words, punctuation marks, and texts are naturally appearing in a sequential order.

Let us introduce RNN in a more formal manner. As aforementioned, it does not just feed-forward, and it even contains network loops other than proceeding from the input layer to the output layer, the network ascends along the layer hierarchy, it also processes data with the time passing measured in $t = 0, 1, \dots$. Each neuron j of a recurrent (hidden) layer d_ℓ has a real-valued vector state $h_{j,t}^{(\ell)}$, which can be regarded as the memory of this neuron unit j . By then, each unit j in each layer ℓ receives two sources of inputs:

1. $\mathbf{h}_t^{(\ell-1)}$: A vector of states from the previous layer $\ell - 1$ at the time step t ;
2. $\mathbf{h}_{t-1}^{(\ell)}$: A vector of states of output from this present layer ℓ but at the previous time step $t - 1$.

11.1 More on Matrix Calculus

In Section 3.1, we have defined a set of notations for the gradient of a function f with respect to a vector or a matrix. However, since the backpropagation through time in RNNs requires a heavy use of chain rules for functions taking values as vector or matrix, for the sake of convenience, we here adopt another set of notations:

1. If $\mathbf{x} \in \mathbb{R}^N$ is a vector variable and $f(\mathbf{x})$ is a real-valued function differentiable in \mathbf{x} , then

$$\nabla_{\mathbf{x}} f(\mathbf{x}) := \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_N} \end{pmatrix} \in \mathbb{R}^N.$$

2. If $\mathbf{x} \in \mathbb{R}^N$ is a vector variable and $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^M$ is a vector-valued function differentiable in \mathbf{x} , then

$$\nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}) := \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_M(\mathbf{x})}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_1(\mathbf{x})}{\partial x_N} & \dots & \frac{\partial f_M(\mathbf{x})}{\partial x_N} \end{pmatrix} = \left(\begin{array}{c|c|c} \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}} & \frac{\partial f_2(\mathbf{x})}{\partial \mathbf{x}} & \dots & \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}} \end{array} \right) \in \mathbb{R}^{N \times M}.$$

3. If $\mathbf{X} \in \mathbb{R}^{N \times M}$ is a matrix variable and $\mathbf{f}(\mathbf{X}) \in \mathbb{R}^K$ is a vector-valued function differentiable in \mathbf{X} , then

$$\nabla_{\mathbf{X}} \mathbf{f}(\mathbf{X}) := \frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} = \begin{pmatrix} \frac{\partial \mathbf{f}(\mathbf{X})}{\partial x_{1,1}} & \dots & \frac{\partial \mathbf{f}(\mathbf{X})}{\partial x_{1,M}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mathbf{f}(\mathbf{X})}{\partial x_{N,1}} & \dots & \frac{\partial \mathbf{f}(\mathbf{X})}{\partial x_{N,M}} \end{pmatrix} \in \mathbb{R}^{N \times M \times K},$$

which is a 3^{rd} order tensor.

In backpropagation of MLP, we only consider the chain rule in scalar form as discussed in Subsubsection 9.1.1.2. In general, the chain rule against the independent variable $x \in \mathbb{R}$ on another scalar function \mathcal{E} that depends x through another function h is

$$\frac{\partial \mathcal{E}}{\partial x} = \frac{\partial \mathcal{E}}{\partial h} \frac{\partial h}{\partial x} = \frac{\partial h}{\partial x} \frac{\partial \mathcal{E}}{\partial h}$$

However, when dealing with a column vector $\mathbf{x} \in \mathbb{R}^N$, the chain rule builds “from the right to the left”, i.e.

$$\frac{\partial \mathcal{E}}{\partial \mathbf{x}} = \frac{\partial \mathcal{E}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{x}}, \quad (11.1.1)$$

the product here is a matrix multiplication. Before we proceed forward, we first introduce three more mathematical operations:

Definition 11.1.1 (Kronecker product \otimes). Given two real matrices $\mathbf{A} = (a_{n,m})_{n=1,m=1}^{N,M}$ of dimensions $N \times M$ and $\mathbf{B} = (b_{k\ell})_{k=1,\ell=1}^{K,L}$ of dimensions $K \times L$, their Kronecker product $\mathbf{A} \otimes \mathbf{B}$ is defined as another matrix of dimensions $(N \cdot K) \times (M \cdot L)$ that

$$\mathbf{A} \otimes \mathbf{B} = \begin{pmatrix} a_{11}\mathbf{B} & \dots & a_{1M}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{N1}\mathbf{B} & \dots & a_{NM}\mathbf{B} \end{pmatrix} \in \mathbb{R}^{(N \cdot K) \times (M \cdot L)},$$

with entries, for $n = 1, \dots, N$, $m = 1, \dots, M$, $k = 1, \dots, K$, and $\ell = 1, \dots, L$,

$$(\mathbf{A} \otimes \mathbf{B})_{K \cdot (n-1) + k, L \cdot (m-1) + \ell} := a_{nm} \cdot b_{k\ell}.$$

Definition 11.1.2 (*vectorization*). Let $\mathbf{X} = (\mathbf{x}_{nm})_{n=1, m=1}^{N, M}$ be a third order tensor in $\mathbb{V}^{N \times M}$, where \mathbf{x}_{nm} takes vector-value in a vector space \mathbb{V} , then the vectorized form $\text{vec}(\mathbf{X})$ of \mathbf{X} is a second order tensor

$$\text{vec}(\mathbf{X}) := (\mathbf{x}_{11}, \dots, \mathbf{x}_{N1}, \mathbf{x}_{12}, \dots, \mathbf{x}_{N2}, \dots, \mathbf{x}_{1M}, \dots, \mathbf{x}_{NM})^\top \in \mathbb{V}^{N \cdot M}.$$

That is to say we append consecutive columns together to form a long column vector. Besides, its inverse operation $\text{vec}^{-1} : \mathbb{V}^{N \cdot M} \rightarrow \mathbb{V}^{N \times M}$ can be defined by splitting the long column of entries into equal length of N of sub-vectors, each of which are then arranged as consecutive columns in order to form a third order tensor, and so

$$\text{vec}^{-1}[\text{vec}(\mathbf{X})] = \mathbf{X} = \begin{pmatrix} \mathbf{x}_{11} & \cdots & \mathbf{x}_{1M} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_{N1} & \cdots & \mathbf{x}_{NM} \end{pmatrix} \in \mathbb{V}^{N \times M}.$$

Lemma 11.1.1. Let $\mathbf{A} \in \mathbb{R}^{K \times N}$, $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_M) \in \mathbb{R}^{N \times M}$, where each $\mathbf{x}_m \in \mathbb{R}^N$ for $m = 1, \dots, M$, and $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_L) \in \mathbb{R}^{M \times L}$ so that $b_\ell \in \mathbb{R}^M$ for $\ell = 1, \dots, L$, then we have

$$\text{vec}(\mathbf{AXB}) = (\mathbf{B}^\top \otimes \mathbf{A})\text{vec}(\mathbf{X}).$$

Proof. From Definition 11.1.2, we have $\mathbb{V} = \mathbb{R}$, we then consider the ℓ^{th} column, for $\ell = 1, \dots, L$, of \mathbf{AXB} :

$$\begin{aligned} (\mathbf{AXB})_{\cdot, \ell} &= \mathbf{AX}\mathbf{b}_\ell = \mathbf{A} \sum_{m=1}^M \mathbf{x}_m b_{m, \ell} = \sum_{m=1}^M (b_{m, \ell} \mathbf{A}) \mathbf{x}_m = \left(\begin{array}{cccc} b_{1, \ell} \mathbf{A} & b_{2, \ell} \mathbf{A} & \cdots & b_{M, \ell} \mathbf{A} \end{array} \right) \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_M \end{pmatrix} \\ &= \left(\left(\begin{array}{cccc} b_{1, \ell} & b_{2, \ell} & \cdots & b_{M, \ell} \end{array} \right) \otimes \mathbf{A} \right) \text{vec}(\mathbf{X}) = (\mathbf{b}_\ell^\top \otimes \mathbf{A}) \text{vec}(\mathbf{X}). \end{aligned}$$

Finally, stacking the columns together vertically, we see that

$$\mathbf{AXB} = \begin{pmatrix} \mathbf{b}_1^\top \otimes \mathbf{A} \\ \mathbf{b}_2^\top \otimes \mathbf{A} \\ \vdots \\ \mathbf{b}_L^\top \otimes \mathbf{A} \end{pmatrix} \text{vec}(\mathbf{X}) = \left[\begin{pmatrix} \mathbf{b}_1^\top \\ \mathbf{b}_2^\top \\ \vdots \\ \mathbf{b}_L^\top \end{pmatrix} \otimes \mathbf{A} \right] \text{vec}(\mathbf{X}) = (\mathbf{B}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{X}),$$

therefore the claim is concluded. \square

Lemma 11.1.2. Let $\mathbf{f}(\mathbf{X}) : \mathbb{R}^{N \times M} \rightarrow \mathbb{R}^K$ be a vector-valued function differentiable in a matrix variable \mathbf{X} , then $(f \circ \text{vec}^{-1})$ is differentiable in $\text{vec}(\mathbf{X})$, and as a long vector in $(\mathbb{R}^K)^{N \times M}$ with an entries in $\mathbb{V} = (\mathbb{R}^K)^\top$,

$$\text{vec} \left(\frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \right) = \frac{\partial \mathbf{f}(\mathbf{X})}{\partial \text{vec}(\mathbf{X})} \in [(\mathbb{R}^K)^\top]^{N \times M} = \mathbb{R}^{N \times M \times K}.$$

Proof. Using Definition 11.1.2, since

$$\frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} = \left(\frac{\partial f_k(\mathbf{X})}{\partial x_{nm}} \right)_{n=1, m=1, k=1}^{N, M, K} = \left(\frac{\partial \mathbf{f}(\mathbf{X})^\top}{\partial x_{nm}} \right)_{n=1, m=1}^{N, M}, \quad \text{where } \frac{\partial \mathbf{f}(\mathbf{X})}{\partial x_{nm}} \in \mathbb{V} = \mathbb{R}^K,$$

we have

$$\text{vec}\left(\frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}}\right) = \left(\begin{array}{c|c|c|c} \frac{\partial \mathbf{f}(\mathbf{X})^\top}{\partial x_{11}} & \dots & \frac{\partial \mathbf{f}(\mathbf{X})^\top}{\partial x_{N1}} & \left| \begin{array}{c|c|c|c} \frac{\partial \mathbf{f}(\mathbf{X})^\top}{\partial x_{12}} & \dots & \frac{\partial \mathbf{f}(\mathbf{X})^\top}{\partial x_{N2}} & \dots \\ \vdots & \ddots & \vdots & \vdots \\ \frac{\partial \mathbf{f}(\mathbf{X})^\top}{\partial x_{1M}} & \dots & \frac{\partial \mathbf{f}(\mathbf{X})^\top}{\partial x_{NM}} & \end{array} \right. \end{array} \right)$$

$$= \frac{\partial \mathbf{f}(\mathbf{X})}{\partial \text{vec}(\mathbf{X})}.$$

□

Remark 11.1.1. To ease the notational complexity, let $\mathbf{A} = (a_{nm})_{n=1,m=1}^{N,M} \in \mathbb{R}^{N \times M}$ and $\mathbf{v} \in \mathbb{R}^K$, in the following we abuse the notation that the third order tensor can be rewritten:

$$\begin{pmatrix} a_{11}\mathbf{v}^\top & \dots & a_{1M}\mathbf{v}^\top \\ \vdots & \ddots & \vdots \\ a_{N1}\mathbf{v}^\top & \dots & a_{NM}\mathbf{v}^\top \end{pmatrix} := \begin{pmatrix} a_{11} & \dots & a_{1M} \\ \vdots & \ddots & \vdots \\ a_{N1} & \dots & a_{NM} \end{pmatrix} \mathbf{v}^\top \in [(\mathbb{R}^K)^\top]^{N \times M} = \mathbb{R}^{N \times M \times K}.$$

Using the afore-defined notations, we here provide some examples that are commonly evaluated in DNN:

1. Note that the derivative of a vector $\mathbf{x} \in \mathbb{R}^N$ with respect to its own is the identity matrix \mathbf{I}_N . Therefore, if we differentiate $\mathbf{x} \odot \mathbf{b}$ with respect to \mathbf{x} , then

$$\frac{\partial \mathbf{x} \odot \mathbf{b}}{\partial \mathbf{x}} = \text{diag}\{\mathbf{b}\} \in \mathbb{R}^{N \times N}, \quad (11.1.2)$$

where $\text{diag}\{\mathbf{b}\}$ is the diagonal matrix of the vector $\mathbf{b} = (b_1, \dots, b_N)^\top \in \mathbb{R}^N$, i.e.

$$\text{diag}\{\mathbf{b}\} = \begin{pmatrix} b_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & b_N \end{pmatrix} \in \mathbb{R}^{N \times N}.$$

Moreover, if we multiply it with a vector $\mathbf{c} \in \mathbb{R}^N$, then

$$\text{diag}\{\mathbf{b}\} \mathbf{c} = \mathbf{b} \odot \mathbf{c}. \quad (11.1.3)$$

2. In the training of MLP, we often evaluate $\mathbf{f}(\mathbf{x}) = \sigma(\mathbf{Ax})$, where $\mathbf{A} \in \mathbb{R}^{M \times N}$ does not involve \mathbf{x} , using (11.1.1), then

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \frac{\partial \mathbf{Ax}}{\partial \mathbf{x}} \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{Ax}} = \mathbf{A}^\top \frac{\partial \sigma(\mathbf{Ax})}{\partial (\mathbf{Ax})},$$

where we recall Property 9 in Section 3.4. Moreover, in MLP, $\sigma_i(\mathbf{Ax})$ and $(\mathbf{Ax})_j$ are often assumed to be independent for $i \neq j$, while $\sigma_i(\mathbf{Ax})$ solely depends on $(\mathbf{Ax})_i$, then we can further reduce the derivative to

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{A}^\top \text{diag}\{\sigma'(\mathbf{Ax})\}.$$

3. To be seen in most RNN problems, consider $\mathbf{f}(\mathbf{X}) = \mathbf{X}\mathbf{a} \in \mathbb{R}^N$, where $\mathbf{a} \in \mathbb{R}^M$ is a fixed vector and $\mathbf{X} \in \mathbb{R}^{N \times M}$ is a matrix variable. We then consider the gradient of $\mathbf{f}(\mathbf{X})$ with respect to a vectorized \mathbf{X} , $\text{vec}(\mathbf{X}) = (x_{11}, \dots, x_{N1}, x_{12}, \dots, x_{N2}, \dots, x_{1M}, \dots, x_{NM})^\top \in \mathbb{R}^{N \cdot M}$. By Lemma 11.1.1, we can rewrite, being a \mathbb{R}^N -vector, $\mathbf{f}(\mathbf{X}) = \mathbf{I}_N \mathbf{X} \mathbf{a} = \text{vec}(\mathbf{I}_N \mathbf{X} \mathbf{a}) = (\mathbf{a}^\top \otimes \mathbf{I}_N) \text{vec}(\mathbf{X})$, then one can use, again, Property 9 in Section 3.4,

$$\frac{\partial \mathbf{f}(\mathbf{X})}{\partial \text{vec}(\mathbf{X})} = \frac{\partial (\mathbf{a}^\top \otimes \mathbf{I}_N) \text{vec}(\mathbf{X})}{\partial \text{vec}(\mathbf{X})} = (\mathbf{a}^\top \otimes \mathbf{I}_N)^\top = \left(a_1 \mathbf{I}_N \mid \dots \mid a_M \mathbf{I}_N \right)^\top \in [(\mathbb{R}^N)^\top]^{(N \cdot M)} = \mathbb{R}^{(N \cdot M) \times N}.$$

By Lemma 11.1.2, we have

$$\frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} = \text{vec}^{-1} \left[\text{vec} \left(\frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \right) \right] = \text{vec}^{-1} \left[\frac{\partial \mathbf{f}(\mathbf{X})}{\partial \text{vec}(\mathbf{X})} \right] = \text{vec}^{-1} (\mathbf{a} \otimes \mathbf{I}_N) \in [(\mathbb{R}^N)^\top]^{N \times M} = \mathbb{R}^{N \times M \times N},$$

which is actually a third order tensor; to further see its effect, let $\mathbf{e}_k = (0, \dots, 0, 1, 0, \dots, 0)^\top \in \mathbb{R}^N$ taking a value of 1 only at the k^{th} position or otherwise zero, then

$$\begin{aligned} \left(\begin{array}{c} a_1 \mathbf{I}_N \\ a_2 \mathbf{I}_N \\ \vdots \\ a_M \mathbf{I}_N \end{array} \right) &\xrightarrow{\text{vec}^{-1}} \left(\begin{array}{cccc} a_1 \mathbf{I}_N & a_2 \mathbf{I}_N & \cdots & a_M \mathbf{I}_N \end{array} \right) = \left(\begin{array}{cccc} a_1 \mathbf{e}_1^\top & a_2 \mathbf{e}_1^\top & \cdots & a_M \mathbf{e}_1^\top \\ a_1 \mathbf{e}_2^\top & a_2 \mathbf{e}_2^\top & \cdots & a_M \mathbf{e}_2^\top \\ \vdots & \vdots & \ddots & \vdots \\ a_1 \mathbf{e}_N^\top & a_2 \mathbf{e}_N^\top & \cdots & a_M \mathbf{e}_N^\top \end{array} \right) \\ &= \sum_{k=1}^N \left(\begin{array}{c|c|c|c} a_1 \mathbf{e}_k & a_2 \mathbf{e}_k & \cdots & a_M \mathbf{e}_k \end{array} \right) \mathbf{e}_k^\top \in \left[(\mathbb{R}^N)^\top \right]^{N \times M}, \end{aligned}$$

defined in Remark 11.1.1, where, for $k = 1, \dots, N$,

$$\left(\begin{array}{c|c|c|c} a_1 \mathbf{e}_k & a_2 \mathbf{e}_k & \cdots & a_M \mathbf{e}_k \end{array} \right) = \left(\begin{array}{cccc} 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ a_1 & a_2 & \cdots & a_M \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{array} \right) \in \mathbb{R}^{N \times M},$$

which only have the k^{th} row with possibly non-zero entries; also see Figure 11.1.1. Moreover, if we further multiply $\nabla_{\mathbf{X}} \mathbf{f}(\mathbf{X})$ by a vector $\mathbf{b} \in \mathbb{R}^N$ on right, then

$$\begin{aligned} \frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \mathbf{b} &= \frac{\partial \mathbf{X} \mathbf{a}}{\partial \mathbf{X}} \mathbf{b} = \text{vec}^{-1}(\mathbf{a} \otimes \mathbf{I}_N) \mathbf{b} = \left(\begin{array}{cccc} a_1 \mathbf{I}_N & a_2 \mathbf{I}_N & \cdots & a_M \mathbf{I}_N \end{array} \right) \mathbf{b} \\ &= \left(\begin{array}{c|c|c|c} a_1 \mathbf{b} & a_2 \mathbf{b} & \cdots & a_M \mathbf{b} \end{array} \right) = \mathbf{b} \mathbf{a}^\top \in \mathbb{R}^{N \times M}. \end{aligned} \quad (11.1.4)$$

First slice coordinates for \mathbf{e}_1^\top				Second slice coordinates for \mathbf{e}_2^\top				M^{th} slice coordinates for \mathbf{e}_N^\top			
a_1	a_2	\cdots	a_M	0	0	\cdots	0	0	0	\cdots	0
0	0	\cdots	0	a_1	a_2	\cdots	a_M	0	0	\cdots	0
\vdots	\vdots	\ddots	\vdots	\vdots	\vdots	\ddots	\vdots	\vdots	\vdots	\ddots	\vdots
0	0	\cdots	0	0	0	\cdots	0	a_1	a_2	\cdots	a_M

Figure 11.1.1: N slices of the first two dimensions in $\mathbb{R}^{N \times M}$ of $\nabla_{\mathbf{X}} \mathbf{f}(\mathbf{X}) = \nabla_{\mathbf{X}} \mathbf{X} \mathbf{a}$.

Example 11.1.1. Using the set of notations above and (11.1.1), the backpropagation of a MLP discussed in Subsubsection 9.1.1.2 can be rewritten as:

$$\begin{aligned} \boldsymbol{\delta}^{(L,t)} &:= \frac{\partial \mathcal{E}^{(t)}}{\partial \mathbf{z}^{(L,t)}} = \frac{\partial \mathbf{a}^{(L,t)}}{\partial \mathbf{z}^{(L,t)}} \frac{\partial \mathcal{E}^{(t)}}{\partial \mathbf{a}^{(L,t)}} = \frac{\partial \mathbf{f}_L(\mathbf{z}^{(L,t)})}{\partial \mathbf{z}^{(L,t)}} \frac{\partial}{\partial \mathbf{a}^{(L,t)}} \left(\frac{1}{2} (\mathbf{y}^{(t)} - \mathbf{a}^{(L,t)})^\top (\mathbf{y}^{(t)} - \mathbf{a}^{(L,t)}) \right) \\ &= \text{diag} \left\{ \mathbf{f}'_L(\mathbf{z}^{(L,t)}) \right\} (\mathbf{a}^{(L,t)} - \mathbf{y}^{(t)}) = \mathbf{f}'_L(\mathbf{z}^{(L,t)}) \odot (\mathbf{a}^{(L,t)} - \mathbf{y}^{(t)}), \end{aligned}$$

the second last equality is directly followed from (11.1.2). Moreover,

$$\begin{aligned}\boldsymbol{\delta}^{(\ell,t)} &:= \frac{\partial \mathcal{E}^{(t)}}{\partial \mathbf{z}^{(\ell,t)}} = \frac{\partial \mathbf{z}^{(\ell+1,t)}}{\partial \mathbf{z}^{(\ell,t)}} \frac{\partial \mathcal{E}^{(t)}}{\partial \mathbf{z}^{(\ell+1,t)}} = \frac{\partial \mathbf{a}^{(\ell,t)}}{\partial \mathbf{z}^{(\ell,t)}} \frac{\partial \mathbf{z}^{(\ell+1,t)}}{\partial \mathbf{a}^{(\ell,t)}} \frac{\partial \mathcal{E}^{(t)}}{\partial \mathbf{z}^{(\ell+1,t)}} \\ &= \text{diag}\{\mathbf{f}'_\ell(\mathbf{z}^{(\ell,t)})\} \left(\mathbf{W}^{(\ell+1,t-1)}\right)^\top \boldsymbol{\delta}^{(\ell+1,t)} = \left\{ \left(\mathbf{W}^{(\ell+1,t-1)}\right)^\top \boldsymbol{\delta}^{(\ell+1,t)} \right\} \odot \mathbf{f}'_\ell(\mathbf{z}^{(\ell,t)}),\end{aligned}$$

where the second last equality follows by (11.1.2) for the first term and for the second term by Property 9 in Section 3.4. Finally, for the weights, by using (11.1.4), we have

$$\frac{\partial \mathcal{E}^{(t)}}{\partial \mathbf{W}^{(\ell,t-1)}} = \frac{\partial \mathbf{z}^{(\ell,t)}}{\partial \mathbf{W}^{(\ell,t-1)}} \frac{\partial \mathcal{E}^{(t)}}{\partial \mathbf{z}^{(\ell,t)}} = \frac{\partial (\mathbf{W}^{(\ell,t-1)} \mathbf{a}^{(\ell-1,t)})}{\partial \mathbf{W}^{(\ell,t-1)}} \boldsymbol{\delta}^{(\ell,t)} = \boldsymbol{\delta}^{(\ell,t)} \left(\mathbf{a}^{(\ell-1,t)}\right)^\top = \boldsymbol{\delta}^{(\ell,t)} \left(\mathbf{f}_{\ell-1}(\mathbf{z}^{(\ell-1,t)})\right)^\top;$$

and then that for biases

$$\frac{\partial \mathcal{E}^{(t)}}{\partial \mathbf{b}^{(\ell,t-1)}} = \frac{\partial \mathbf{z}^{(\ell,t)}}{\partial \mathbf{b}^{(\ell,t-1)}} \frac{\partial \mathcal{E}^{(t)}}{\partial \mathbf{z}^{(\ell,t)}} = \frac{\partial \mathbf{b}^{(\ell,t-1)}}{\partial \mathbf{b}^{(\ell,t-1)}} \boldsymbol{\delta}^{(\ell,t)} = \mathbf{I}_{d_\ell} \boldsymbol{\delta}^{(\ell,t)} = \boldsymbol{\delta}^{(\ell,t)}.$$

11.2 Process in RNN

As aforementioned, RNN takes a sequential data matrix \mathbf{X} as input, in which each datum is essentially a sequence of vectors $\mathbf{x}_t \in \mathbb{R}^d$, for $t = 1, \dots, D$ and some $d \in \mathbb{Z}^+$, and we say that the data input is of a length D :

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{t-1}, \mathbf{x}_t, \mathbf{x}_{t+1}, \dots, \mathbf{x}_D).$$

For instance, if the input sample \mathbf{X} is a text sentence, then the feature vector \mathbf{x}_t , for each $t = 1, \dots, D$, represents a word in a sentence at time t . For instance, Table 11.2.1 shows a sequence of input feature vectors of length 8 in a movie review for sentiment analysis, in which a positive review indicating a favourable consideration has a label of $y = 1$; otherwise, $y = 0$, implies an inferior performance.

Data	\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_3	\mathbf{x}_4	\mathbf{x}_5	\mathbf{x}_6	\mathbf{x}_7	\mathbf{x}_8	\mathbf{y}
1	The	first	half	was	very	boring	.		0
2	Great	performance	by	all	the	leading	actors	.	1
3	The	visual	effects	were	stunning	.			1
4	The	movie	was	a	waste	of	time	.	0

Table 11.2.1: The inputs and labels of an RNN in a movie review for sentiment analysis.

11.2.1 Forward Propagation of a Temporal Dataset

The forward propagation of a L -layer RNN takes the form, for $t = 1, 2, \dots, D$,

$$(Cell\ output) \quad \mathbf{h}_t^{(\ell)} = \mathbf{f}_\ell(\mathbf{W}_h^{(\ell)} \mathbf{h}_t^{(\ell-1)} + \mathbf{V}_h^{(\ell)} \mathbf{h}_{t-1}^{(\ell)} + \mathbf{b}_h^{(\ell)}), \quad \text{for } \ell = 1, \dots, L-1; \quad (11.2.1a)$$

$$(Final\ output) \quad \mathbf{o}_t = \sigma(\mathbf{W}_o \mathbf{h}_t^{(L-1)} + \mathbf{b}_o). \quad (11.2.1b)$$

For any single layer ℓ , every neuron or now called *memory cell* uses the same set of coefficients $\mathbf{W}_h^{(\ell)}, \mathbf{V}_h^{(\ell)}$, and $\mathbf{b}_h^{(\ell)}$ being independent of time t , it means that the model is temporal homogeneous. For example, Figure 11.2.1 illustrates the first three timestep of a 2-layer RNN with the activation function \mathbf{f}_ℓ , for $\ell = 1, 2$, being the hyperbolic tangent function:

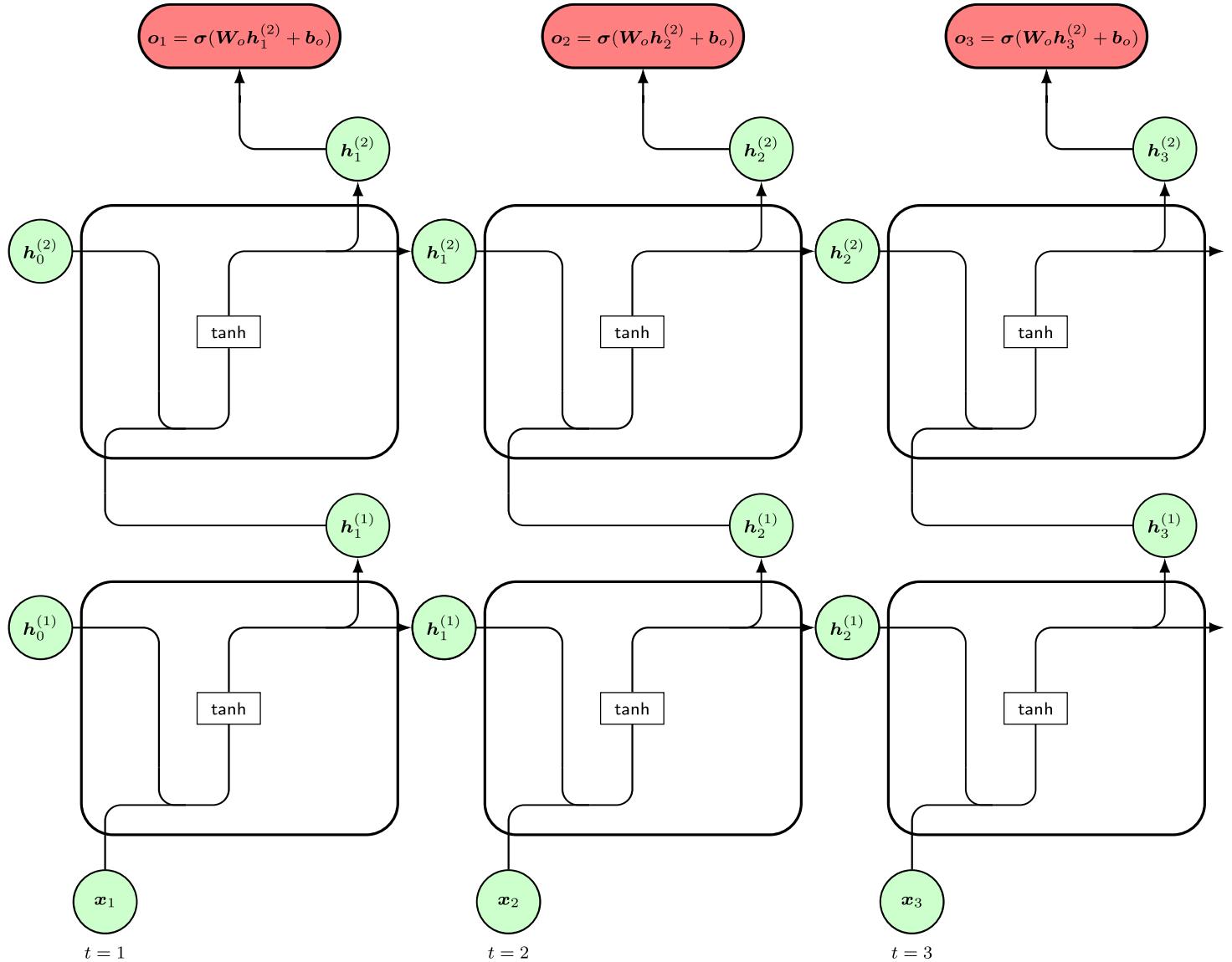


Figure 11.2.1: The first three timestep of a two-layer RNN; also see [4].

The dimension of $\mathbf{W}_o^{(\ell)}$ is chosen such that the multiplicative product of the matrix $\mathbf{W}_o^{(\ell)}$ and the vector $\mathbf{h}_t^{(\ell)}$ has the same dimension as that of the bias vector $\mathbf{b}_o^{(\ell)}$; in return, their choices also depend on the very dimension of the output label \mathbf{y} in the training dataset. In a typical MLP, the weights and biases are updated through the backpropagation of the training dataset, as introduced in Subsubsection 9.1.1.2; but for this temporal dependent RNN, a modified version of the usual backpropagation is now considered to update $\mathbf{W}_h^{(\ell)}, \mathbf{V}_h^{(\ell)}, \mathbf{b}_h^{(\ell)}, \mathbf{W}_o^{(\ell)}$, and $\mathbf{b}_o^{(\ell)}$, which is called *backpropagation through time*; see Subsection 11.2.3 for a detailed discussion. Lastly, the initial states $\mathbf{h}_0^{(\ell)}$ can be fixed parameters, *i.e.* initializing $\mathbf{h}_0^{(\ell)}$ with random numbers, or learnable parameter, *i.e.* initializing $\mathbf{h}_0^{(\ell,0)}$ with random numbers and then updating $\mathbf{h}_0^{(\ell,k)}$ for each iteration k with backpropagation through time. For simplicity, to avoid unnecessary technicalities, we shall only consider fixed initial states $\mathbf{h}_0^{(\ell)}$ in the rest of our discussions.

11.2.2 An Illustrative Example of RNN

Figure 11.2.2 considers a particular timestep t of the 2-layer RNN shown in Figure 11.2.1:

1. The first (leftmost) layer receives as input a feature vector $\mathbf{x}_t = (x_{1,t}, x_{2,t})^\top \in \mathbb{R}^2$ and the output of the previous timestep from the same layer $\mathbf{h}_{t-1}^{(\ell)} = (h_{1,t-1}^{(\ell)}, h_{2,t-1}^{(\ell)}) \in \mathbb{R}^2$; and
2. The second layer receives the output of the first layer as input $\mathbf{h}_t^{(\ell-1)} = (h_{1,t}^{(\ell-1)}, h_{2,t}^{(\ell-1)})^\top \in \mathbb{R}^2$.

Therefore, for $d = 2$, we have, for $\ell = 1, 2$,

$$\mathbf{W}_h^{(\ell)} = \begin{pmatrix} \omega_{h,1}^{(\ell)} & \omega_{h,2}^{(\ell)} \end{pmatrix}^\top \quad \text{and} \quad \mathbf{b}_h^{(\ell)} = (b_{h,1}^{(\ell)}, b_{h,2}^{(\ell)})^\top.$$

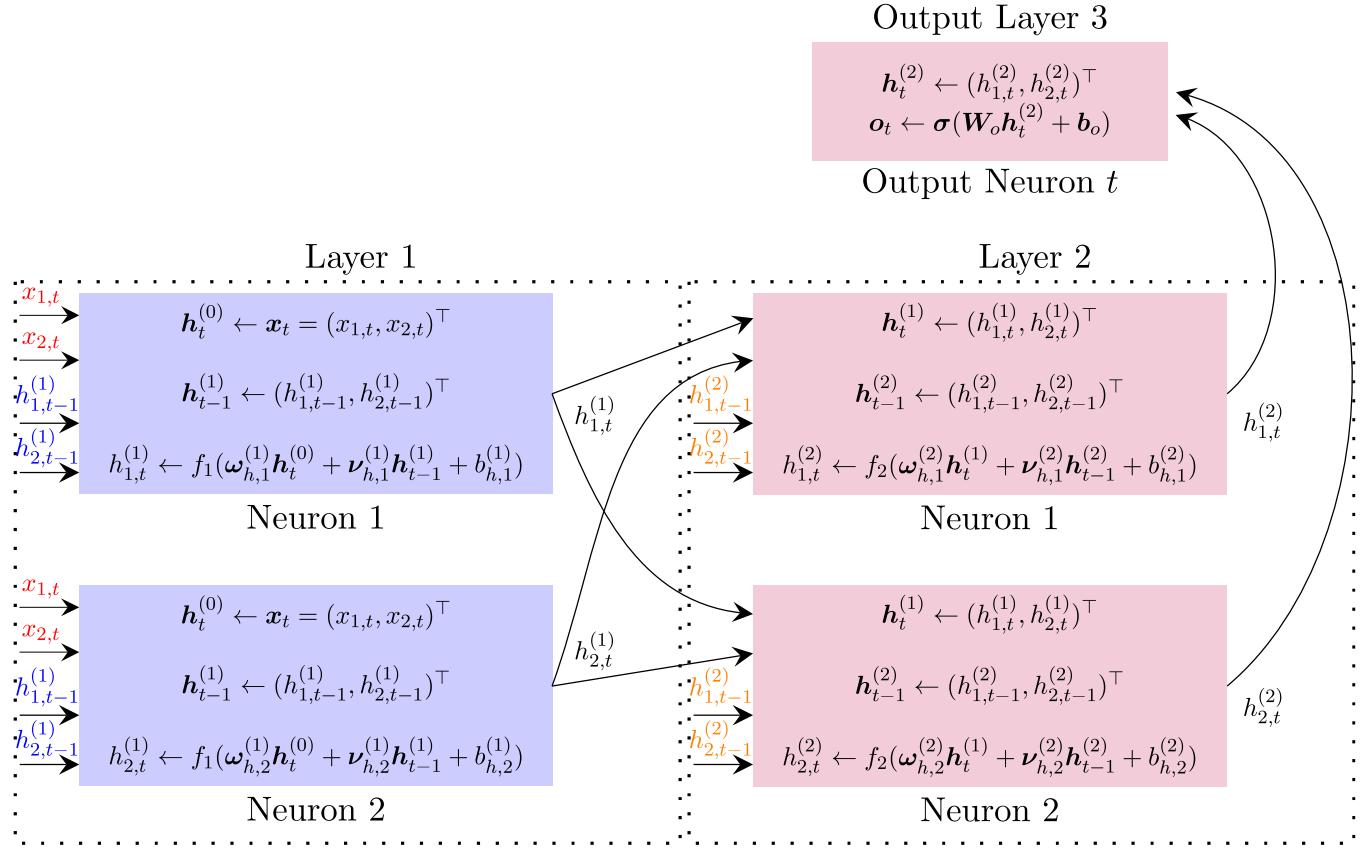


Figure 11.2.2: A folded RNN: The input feature vector is two-dimensional, each hidden layer has two neuron units, inspired by [1].

For RNN, to update the state $\mathbf{h}_{j,t}^{(\ell)}$ at each timestep $t = 1, \dots, D$ in each neuron unit j of the layer ℓ :

1. calculate a linear combination of the inputs $\mathbf{h}_t^{(\ell-1)}$ from the previous layer and $\mathbf{h}_{t-1}^{(\ell)}$ from the same layer but the last timestep with the respective parameter vectors $\omega_{h,j}^{(\ell)}$, and $\nu_{h,j}^{(\ell)}$, and the bias $b_{h,j}^{(\ell)}$, in other words, $(\omega_{h,j}^{(\ell)})^\top \mathbf{h}_t^{(\ell-1)} + (\nu_{h,j}^{(\ell)})^\top \mathbf{h}_{t-1}^{(\ell)} + b_{h,j}^{(\ell)}$; then
2. apply an activation function f_1 to obtain the current output $\mathbf{h}_{j,t}^{(\ell)}$ of this neuron unit j , and so $\mathbf{h}_{j,t}^{(\ell)} = f_1[(\omega_{h,j}^{(\ell)})^\top \mathbf{h}_t^{(\ell-1)} + (\nu_{h,j}^{(\ell)})^\top \mathbf{h}_{t-1}^{(\ell)} + b_{h,j}^{(\ell)}]$ with $\mathbf{h}_t^{(0)} = \mathbf{x}_t$, for $t = 1, \dots, D$. Some common choices for f_1 are sigmoid type activation functions, such as the commonly used hyperbolic tangent \tanh ; while ReLU is not so prevalent in RNNs.

Finally, in the output layer $L = 3$, we shall apply another activation function σ to the output of all neurons that returns a different vector \mathbf{o}_t of the same input dimension of the arguments, i.e. $\mathbf{o}_t = \sigma(\mathbf{W}_o \mathbf{h}_t^{(2)} + \mathbf{b}_o)$.

For instance, a typical choice for σ in a classification problem (with d_L classes) is the softmax function:

$$\sigma(\mathbf{z}^{(L)}) = (\sigma_1(\mathbf{z}^{(L)}), \dots, \sigma_{d_L}(\mathbf{z}^{(L)}))^\top, \quad \text{where } \sigma_i(\mathbf{z}^{(L)}) = \frac{\exp\{z_i^{(L)}\}}{\sum_{j=1}^{d_L} \exp\{z_j^{(L)}\}},$$

where d_L is the total number of neuron units in this layer L .

Figure 11.2.3 illustrates an unfolded 2-layer RNN model shown in Figure 11.2.1. We shall use this model to illustrate more about the feed forward procedure.

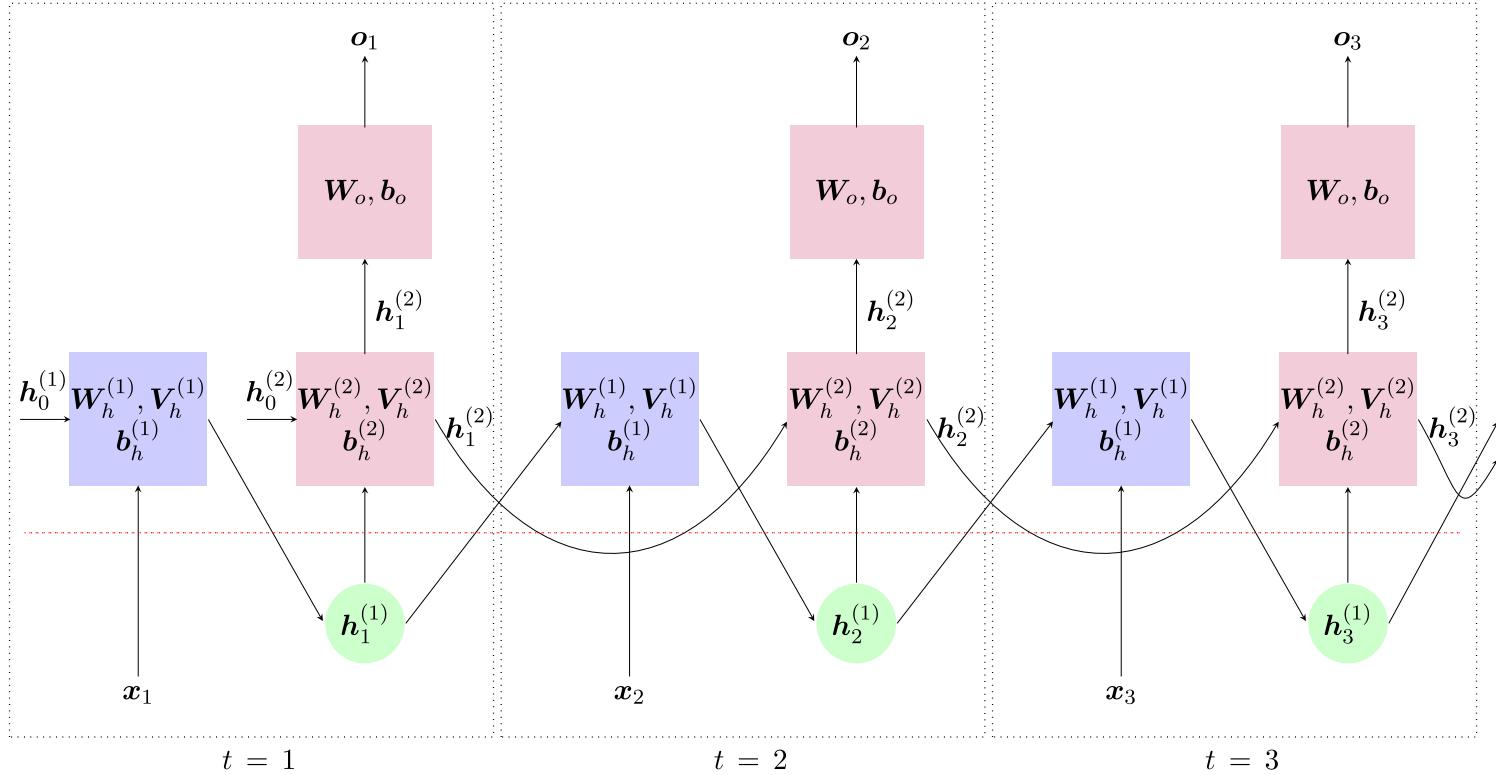


Figure 11.2.3: Unfolding the 2-layer RNN shown in Figure 11.2.1 with $t = 1, 2, 3$, inspired by [1].

11.2.3 Backpropagation through Time

For the sake of convenience, we only illustrate the main idea through one example with MSE as the loss function, and we only consider the calibration for a 1-layer RNN, see Figure 11.2.4. Therefore, the individual temporal error function \mathcal{E}_t at timestep t is simply,

$$\mathcal{E}_t = \frac{1}{2}(\mathbf{y}_t - \mathbf{o}_t)^\top (\mathbf{y}_t - \mathbf{o}_t),$$

where \mathbf{y}_t is the label at time t , \mathbf{o}_t is the predicted output. The total loss function is then given by:

$$\mathcal{E} = \sum_{t=1}^D \mathcal{E}_t.$$

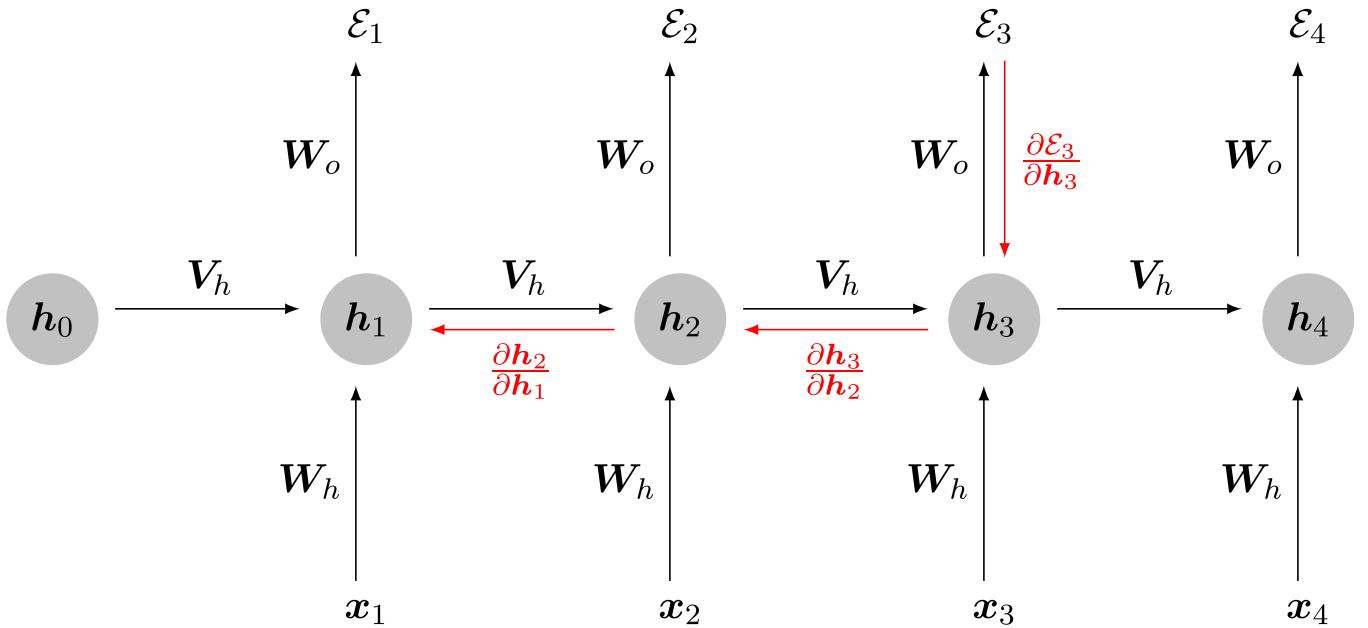


Figure 11.2.4: 1-layer RNN with length $D = 4$: Backpropagation path for $\nabla_{\mathbf{W}_h} \mathcal{E}_3$, $\nabla_{\mathbf{V}_h} \mathcal{E}_3$, and $\nabla_{\mathbf{b}_h} \mathcal{E}_3$.

Recall the notations adopted in (11.2.1a) and (11.2.1b):

$$(Cell\ output) \quad \mathbf{h}_t = \mathbf{f}(\mathbf{W}_h \mathbf{x}_t + \mathbf{V}_h \mathbf{h}_{t-1} + \mathbf{b}_h); \quad (11.2.2a)$$

$$(Final\ output) \quad \mathbf{o}_t = \sigma(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o). \quad (11.2.2b)$$

By the law of total derivative, the gradient for parameters \mathbf{W}_o is:

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \mathbf{W}_o} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{W}_o} = \sum_{t=1}^D \frac{\partial (\mathbf{W}_o \mathbf{h}_t)}{\partial \mathbf{W}_o} \frac{\partial \mathbf{o}_t}{\partial (\mathbf{W}_o \mathbf{h}_t)} \frac{\partial \mathcal{E}_t}{\partial \mathbf{o}_t} = \sum_{t=1}^D \frac{\partial (\mathbf{W}_o \mathbf{h}_t)}{\partial \mathbf{W}_o} \text{diag} \{ \sigma'(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o) \} (\mathbf{o}_t - \mathbf{y}_t) \\ &= \sum_{t=1}^D \{ (\mathbf{o}_t - \mathbf{y}_t) \odot \sigma'(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o) \} \mathbf{h}_t^\top, \end{aligned} \quad (11.2.3)$$

by using (11.1.4); and that of \mathbf{b}_o is

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \mathbf{b}_o} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{b}_o} = \sum_{t=1}^D \frac{\partial \mathbf{o}_t}{\partial \mathbf{b}_o} \frac{\partial \mathcal{E}_t}{\partial \mathbf{o}_t} = \sum_{t=1}^D \text{diag} \{ \sigma'(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o) \} (\mathbf{o}_t - \mathbf{y}_t) \\ &= \sum_{t=1}^D (\mathbf{o}_t - \mathbf{y}_t) \odot \sigma'(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o). \end{aligned} \quad (11.2.4)$$

On the other hand, the gradients of \mathcal{E} with respect to the parameters \mathbf{W}_h , \mathbf{V}_h , and \mathbf{b}_h are more complicated

to derive, because \mathbf{h}_t also depends on the output of the same hidden layer, but at the previous timestep at $t - 1$, *i.e.* \mathbf{h}_{t-1} , while \mathbf{h}_{t-1} itself is also a function of the layer's output, but again even at an earlier timestep $t - 2$, *i.e.* \mathbf{h}_{t-2} , same arguments goes until it reaches the first hidden timestep input \mathbf{h}_0 and the original feature inputs \mathbf{x} 's. And this is precisely what we meant by propagating the error rate back through time.

For the sake of computational convenience, denote the error rate and the weighted sum by, respectively:

$$\delta_{t,\tau} := \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_\tau} \quad \text{and} \quad \mathbf{z}_\tau^h = \mathbf{W}_h \mathbf{x}_\tau + \mathbf{V}_h \mathbf{h}_{\tau-1} + \mathbf{b}_h, \quad \text{for } \tau = 1, \dots, t.$$

Note that directly from (11.2.1b) and Property 9 in Section 3.4,

$$\begin{aligned} \delta_{t,t} &:= \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_t} = \frac{\partial \mathbf{o}_t}{\partial \mathbf{h}_t} \frac{\partial \mathcal{E}_t}{\partial \mathbf{o}_t} = \frac{\partial (\mathbf{W}_o \mathbf{h}_t)}{\partial \mathbf{h}_t} \frac{\partial \mathbf{o}_t}{\partial (\mathbf{W}_o \mathbf{h}_t)} \frac{\partial \mathcal{E}_t}{\partial \mathbf{o}_t} = \mathbf{W}_o^\top \text{diag}\{\sigma'(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o)\} (\mathbf{o}_t - \mathbf{y}_t) \\ &= \mathbf{W}_o^\top \{(\mathbf{o}_t - \mathbf{y}_t) \odot \sigma'(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o)\}. \end{aligned} \quad (11.2.5)$$

The next step is to backpropagate the error rate to the previous timestep $t - 1$:

$$\delta_{t,t-1} := \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_{t-1}} = \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_t} = \frac{\partial \mathbf{z}_t^h}{\partial \mathbf{h}_{t-1}} \frac{\partial \mathbf{h}_t}{\partial \mathbf{z}_t^h} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_t} = \mathbf{V}_h^\top \text{diag}\{\mathbf{f}'(\mathbf{z}_t^h)\} \delta_{t,t} = \mathbf{V}_h^\top \left\{ \delta_{t,t} \odot \mathbf{f}'(\mathbf{z}_t^h) \right\}. \quad (11.2.6)$$

Similarly, inductively, for $\tau = 1, \dots, t - 1$, we obtain

$$\delta_{t,\tau} = \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_\tau} = \frac{\partial \mathbf{z}_{\tau+1}^h}{\partial \mathbf{h}_\tau} \frac{\partial \mathbf{h}_{\tau+1}}{\partial \mathbf{z}_{\tau+1}^h} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_{\tau+1}} = \mathbf{V}_h^\top \left\{ \delta_{t,\tau+1} \odot \mathbf{f}'(\mathbf{z}_{\tau+1}^h) \right\}. \quad (11.2.7)$$

Figure 11.2.4 illustrates the path of the backpropagation through time of $\nabla_{\mathbf{W}_h} \mathcal{E}_3$, $\nabla_{\mathbf{V}_h} \mathcal{E}_3$, and $\nabla_{\mathbf{b}_h} \mathcal{E}_3$, indicated in red, of a 1-layer RNN; in particular, $\nabla_{\mathbf{W}_h} \mathcal{E}_3$ is computed as:

$$\begin{aligned} \frac{\partial \mathcal{E}_3}{\partial \mathbf{W}_h} &= \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_h} \frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathcal{E}_3}{\partial \mathbf{h}_3} + \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}_h} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathcal{E}_3}{\partial \mathbf{h}_3} + \frac{\partial \mathbf{h}_3}{\partial \mathbf{W}_h} \frac{\partial \mathcal{E}_3}{\partial \mathbf{h}_3} = \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_h} \frac{\partial \mathcal{E}_3}{\partial \mathbf{h}_1} + \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}_h} \frac{\partial \mathcal{E}_3}{\partial \mathbf{h}_2} + \frac{\partial \mathbf{h}_3}{\partial \mathbf{W}_h} \frac{\partial \mathcal{E}_3}{\partial \mathbf{h}_3} \\ &= \sum_{\tau=1}^3 \frac{\partial \mathbf{h}_\tau}{\partial \mathbf{W}_h} \frac{\partial \mathcal{E}_3}{\partial \mathbf{h}_\tau} = \sum_{\tau=1}^3 \frac{\partial \mathbf{z}_\tau^h}{\partial \mathbf{W}_h} \frac{\partial \mathbf{h}_\tau}{\partial \mathbf{z}_\tau^h} \frac{\partial \mathcal{E}_3}{\partial \mathbf{h}_\tau}. \end{aligned} \quad (11.2.8)$$

Therefore, using similar calculations of (11.2.3), the gradient of \mathcal{E} with respect to the parameters \mathbf{W}_h , \mathbf{V}_h , and \mathbf{b}_h are, respectively:

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \mathbf{W}_h} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{W}_h} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{z}_\tau^h}{\partial \mathbf{W}_h} \frac{\partial \mathbf{h}_\tau}{\partial \mathbf{z}_\tau^h} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial (\mathbf{W}_h \mathbf{x}_\tau)}{\partial \mathbf{W}_h} \text{diag}\{\mathbf{f}'(\mathbf{z}_\tau^h)\} \delta_{t,\tau} \\ &= \sum_{t=1}^D \sum_{\tau=1}^t \{\mathbf{f}'(\mathbf{z}_\tau^h) \odot \delta_{t,\tau}\} \mathbf{x}_\tau^\top, \end{aligned} \quad (11.2.9)$$

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \mathbf{V}_h} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{V}_h} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{z}_\tau^h}{\partial \mathbf{V}_h} \frac{\partial \mathbf{h}_\tau}{\partial \mathbf{z}_\tau^h} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial (\mathbf{V}_h \mathbf{h}_{\tau-1})}{\partial \mathbf{V}_h} \text{diag}\{\mathbf{f}'(\mathbf{z}_\tau^h)\} \delta_{t,\tau} \\ &= \sum_{t=1}^D \sum_{\tau=1}^t \{\mathbf{f}'(\mathbf{z}_\tau^h) \odot \delta_{t,\tau}\} \mathbf{h}_{\tau-1}^\top, \end{aligned} \quad (11.2.10)$$

$$\frac{\partial \mathcal{E}}{\partial \mathbf{b}_h} = \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{b}_h} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{z}_\tau^h}{\partial \mathbf{b}_h} \frac{\partial \mathbf{h}_\tau}{\partial \mathbf{z}_\tau^h} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \text{diag}\{\mathbf{f}'(\mathbf{z}_\tau^h)\} \delta_{t,\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \mathbf{f}'(\mathbf{z}_\tau^h) \odot \delta_{t,\tau}. \quad (11.2.11)$$

In multilayer RNN, the backpropagations happen through both layers and time. One should pay attention that each $\mathbf{h}_t^{(\ell)}$ is a function of both $\mathbf{h}_{t-1}^{(\ell)}$ and $\mathbf{h}_t^{(\ell-1)}$, while $\mathbf{h}_{t-1}^{(\ell)}$ and $\mathbf{h}_t^{(\ell-1)}$ depend on the pairs $(\mathbf{h}_{t-2}^{(\ell)}, \mathbf{h}_{t-1}^{(\ell-1)})$ and $(\mathbf{h}_{t-1}^{(\ell-1)}, \mathbf{h}_t^{(\ell-2)})$, respectively, the same dependence structure will be carried on until it reaches to the utmost left and floor boundaries; for instance, as depicted in Figure 11.2.5, the interdependence as a function of \mathcal{E} of a 2-layer RNN with $D = 4$, the backprobagation mechanism on $\nabla_{\mathbf{W}_h^{(1)}} \mathcal{E}_3$, $\nabla_{\mathbf{V}_h^{(1)}} \mathcal{E}_3$, and $\nabla_{\mathbf{b}_h^{(1)}} \mathcal{E}_3$ are

illustrated.

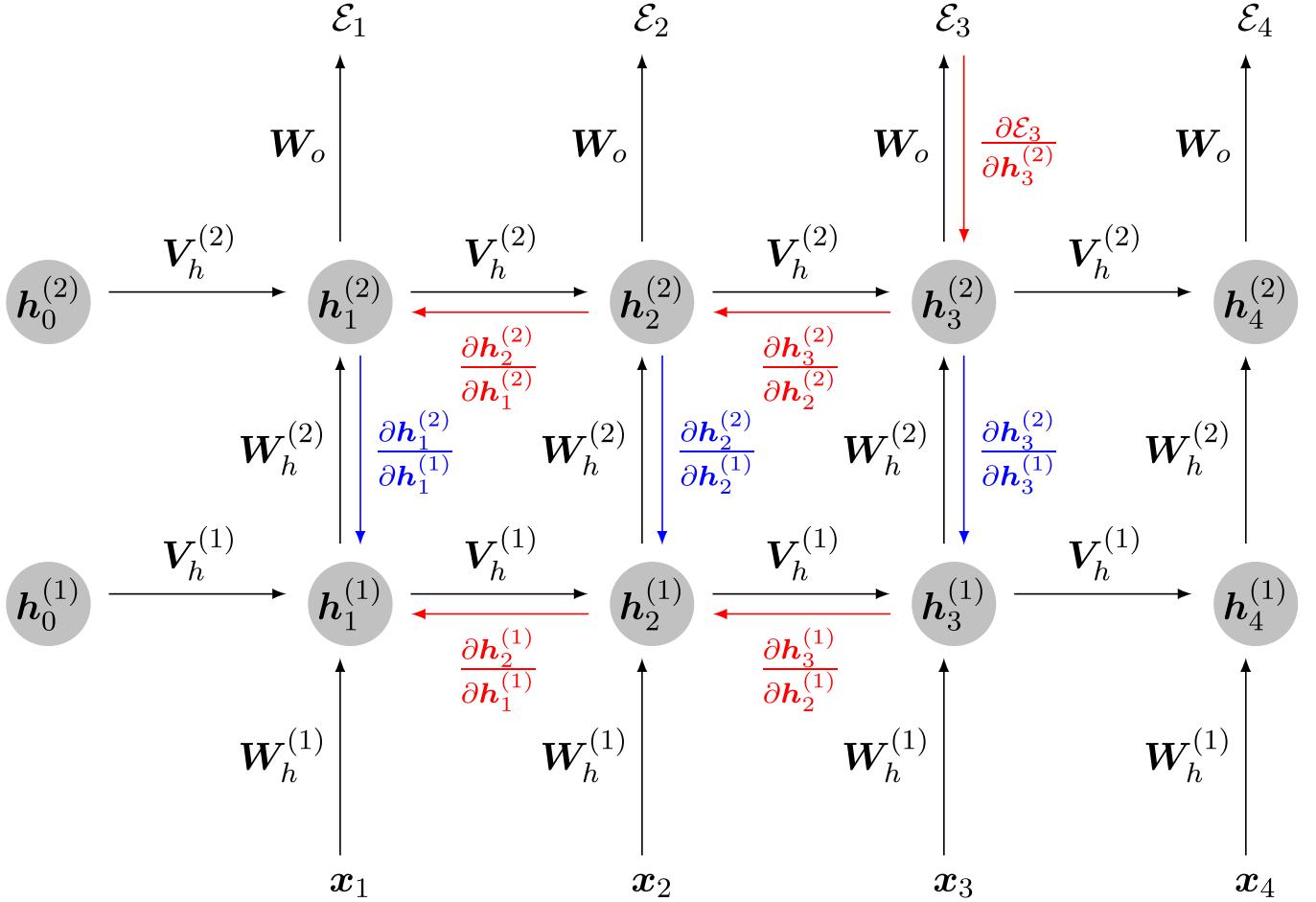


Figure 11.2.5: A 2-layer RNN with length $D = 4$: Backpropagation path of $\nabla_{\mathbf{W}_h^{(1)}} \mathcal{E}_3$, $\nabla_{\mathbf{V}_h^{(1)}} \mathcal{E}_3$, and $\nabla_{\mathbf{b}_h^{(1)}} \mathcal{E}_3$.

From (11.2.5), we know that the gradient of the temporal error \mathcal{E}_t with respect to $\mathbf{h}_t^{(2)}$ and $\mathbf{h}_{\tau}^{(2)}$, for $\tau = 1, \dots, t-1$, are, respectively:

$$\delta_{t,t}^{(2)} = \mathbf{W}_o^\top \left\{ (\mathbf{o}_t - \mathbf{y}_t) \odot \sigma'(\mathbf{W}_o \mathbf{h}_t^{(2)} + \mathbf{b}_o) \right\} \quad \text{and} \quad \delta_{t,\tau}^{(2)} = (\mathbf{V}_h^{(2)})^\top \left\{ \delta_{t,\tau+1}^{(2)} \odot f'(z_{\tau+1}^{h,(2)}) \right\}.$$

We now consider the temporal error \mathcal{E}_t with respect to $\mathbf{h}_t^{(1)}$:

$$\begin{aligned} \delta_{t,t}^{(1)} &= \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_t^{(1)}} = \frac{\partial \mathbf{h}_t^{(2)}}{\partial \mathbf{h}_t^{(1)}} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_t^{(2)}} = \frac{\partial z_t^{h,(2)}}{\partial \mathbf{h}_t^{(1)}} \frac{\partial \mathbf{h}_t^{(2)}}{\partial z_t^{h,(2)}} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_t^{(2)}} \\ &= \frac{\partial (\mathbf{W}_h^{(2)} \mathbf{h}_t^{(1)})}{\partial \mathbf{h}_t^{(1)}} \text{diag}\{f'_2(z_t^{h,(2)})\} \delta_{t,t}^{(2)} = (\mathbf{W}_h^{(2)})^\top \left\{ f'_2(z_t^{h,(2)}) \odot \delta_{t,t}^{(2)} \right\}. \end{aligned}$$

However, for $\delta_{t,t-1}^{(1)}$, we first note in Figure 11.2.5 that there are two paths, one red path and one blue path, leading to the same neuron on the previous timestep $t-1$; to this end, we shall first consider, for $t = 1, \dots, D$,

$$\frac{\partial \mathbf{h}_t^{(2)}}{\partial \mathbf{h}_t^{(1)}} = \frac{\partial z_t^{h,(2)}}{\partial \mathbf{h}_t^{(1)}} \frac{\partial \mathbf{h}_t^{(2)}}{\partial z_t^{h,(2)}} = \frac{\partial (\mathbf{W}_h^{(2)} \mathbf{h}_{t-1}^{(1)})}{\partial \mathbf{h}_{t-1}^{(1)}} \text{diag}\{f'_2(z_t^{h,(2)})\} = (\mathbf{W}_h^{(2)})^\top \text{diag}\{f'_2(z_t^{h,(2)})\},$$

using Property 9 in Section 3.4. Therefore, using which and the intermediate step in (11.2.6), we have

$$\begin{aligned}\boldsymbol{\delta}_{t,t-1}^{(1)} &= \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_{t-1}^{(1)}} = \frac{\partial \mathbf{h}_{t-1}^{(2)}}{\partial \mathbf{h}_{t-1}^{(1)}} \frac{\partial \mathbf{h}_t^{(2)}}{\partial \mathbf{h}_{t-1}^{(2)}} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_t^{(2)}} + \frac{\partial \mathbf{h}_t^{(1)}}{\partial \mathbf{h}_{t-1}^{(1)}} \frac{\partial \mathbf{h}_t^{(2)}}{\partial \mathbf{h}_t^{(1)}} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_t^{(2)}} \\ &= (\mathbf{W}_h^{(2)})^\top \text{diag} \left\{ \mathbf{f}'_2(\mathbf{z}_{t-1}^{h,(2)}) \right\} (\mathbf{V}_h^{(2)})^\top \text{diag} \left\{ \mathbf{f}'_2(\mathbf{z}_t^{h,(2)}) \right\} \boldsymbol{\delta}_{t,t}^{(2)} \\ &\quad + (\mathbf{V}_h^{(1)})^\top \text{diag} \left\{ \mathbf{f}'_1(\mathbf{z}_t^{h,(1)}) \right\} (\mathbf{W}_h^{(2)})^\top \text{diag} \left\{ \mathbf{f}'_2(\mathbf{z}_t^{h,(2)}) \right\} \boldsymbol{\delta}_{t,t}^{(2)}.\end{aligned}$$

For that of the timestep $t - 2$, the computation becomes tedious, we observe in Figure 11.2.5 that there are three paths. In general, for computing the gradient of the temporal error \mathcal{E}_t with respect to $\mathbf{h}_\tau^{(1)}$, for $\tau = 1, \dots, t$, in a 2-layer RNN has $t - \tau + 1$ number of paths, and we shall end the computation here.

11.3 Problems with RNN

- Vanishing /Exploding Gradient:** From Figure 11.2.3, the same set of parameters will be used for calibration for each t of length D time series data, if this length of the input sequence D increases, as explained in last Section 11.2, backpropagation has to “unfold” the network through time. For instance, from (11.2.7), the error rate has to backpropagate through time from timestep D to the timestep 1, namely:

$$\boldsymbol{\delta}_{D,1} = \mathbf{V}_h^\top \left\{ \boldsymbol{\delta}_{D,1+1} \odot \mathbf{f}'(\mathbf{z}_{1+1}^h) \right\} = \underbrace{\mathbf{V}_h^\top \left\{ \dots \left\{ \mathbf{V}_h^\top \left\{ \boldsymbol{\delta}_{D,D} \odot \mathbf{f}'(\mathbf{z}_D^h) \right\} \odot \dots \right\} \odot \mathbf{f}'(\mathbf{z}_2^h) \right\}}_{D-1 \text{ entries}}, \quad (11.3.1)$$

where $\boldsymbol{\delta}_{D,D}$ is given in (11.2.5) with $t = D$; that means the error rate at the timestep 1 is the product of $D - 1$ matrices \mathbf{V}_h^\top . For if \mathbf{V}_h^\top is small and if D is also very large, then this induces a vanishing gradient problem that much affects the parameter update. On the other hand, if \mathbf{V}_h^\top is large, then this will lead to an exploding gradient problem. Note that the problem still exist no matter what the activation function is chosen. On the other hand, as mentioned before in Section 9.2, *tanh* and *softmax* (*sigmoid*) suffer from the vanishing gradient problem in MLP, and the same problem occurs in RNN that the error rate in (11.3.1) is also a product of $\mathbf{f}'(\mathbf{z}_D^h) \odot \dots \odot \mathbf{f}'(\mathbf{z}_2^h)$. Note that, with f being the sigmoid (resp. *tanh*) activation function, we have $f'(\mathbf{z}_t^h) \leq 1/4$ (resp. $f'(\mathbf{z}_t^h) \leq 1$), for $t = 1, \dots, D$. Therefore, the product of $\mathbf{f}'(\mathbf{z}_D^h) \odot \dots \odot \mathbf{f}'(\mathbf{z}_2^h)$ is capped at $(1/4)^{D-1}$ (resp. 1) and this may trigger the vanishing gradient problem.

- Long-term Dependence:** As the length of the time series data grows, the feature vectors from the beginning of the sequence tend to be “forgotten” in the sense that the state of each unit, which serves as network’s memory, is more notably affected by the feature vectors that read recently. Consider the two individual temporal error, \mathcal{E}_1 and \mathcal{E}_D , for the hidden state of \mathbf{h} from timestep 1 and timestep D , respectively, and they are

$$\mathbf{h}_1 = \mathbf{f}(\mathbf{W}_h \mathbf{x}_0 + \mathbf{V}_h \mathbf{h}_0 + \mathbf{b}_h) \quad \text{and} \quad \mathbf{h}_D = \mathbf{f}(\mathbf{W}_h \mathbf{x}_D + \mathbf{V}_h \mathbf{h}_{D-1} + \mathbf{b}_h),$$

while \mathbf{h}_{D-1} is a function of both previous timestep output \mathbf{h}_{D-2} and \mathbf{x}_{D-1} , which in turn is a function of another previous timestep objects and so forth until \mathbf{h}_0 and \mathbf{x}_1 . Hence, as D is large, the impact of \mathbf{x}_1 on \mathbf{h}_D is much smaller than its impact on \mathbf{h}_1 , that means being forgotten; indeed, for 1-layer

RNN, mathematically,

$$\frac{\partial \mathbf{h}_D}{\partial \mathbf{h}_1} = \frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \cdots \frac{\partial \mathbf{h}_{D-1}}{\partial \mathbf{h}_{D-2}} \frac{\partial \mathbf{h}_D}{\partial \mathbf{h}_{D-1}}.$$

Since each

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \frac{\partial \mathbf{z}_t^h}{\partial \mathbf{h}_{t-1}} \frac{\partial \mathbf{h}_t}{\partial \mathbf{z}_t^h} = \mathbf{V}_h^\top \text{diag} \left\{ \mathbf{f}'(\mathbf{z}_t^h) \right\},$$

then

$$\frac{\partial \mathbf{h}_D}{\partial \mathbf{h}_1} = \mathbf{V}_h^\top \text{diag} \left\{ \mathbf{f}'(\mathbf{z}_2^h) \right\} \frac{\partial \mathbf{h}_D}{\partial \mathbf{h}_2} = \mathbf{V}_h^\top \text{diag} \left\{ \mathbf{f}'(\mathbf{z}_2^h) \right\} \mathbf{V}_h^\top \text{diag} \left\{ \mathbf{f}'(\mathbf{z}_3^h) \right\} \cdots \mathbf{V}_h^\top \text{diag} \left\{ \mathbf{f}'(\mathbf{z}_D^h) \right\}.$$

when one applies RNNs to text or speech recognition, the cause-effect link between distant words in a long sentence can be lost. One possible solution to this short memory problem is to use **gated RNNs** or **LSTM** to be discussed in the next two sections.

11.4 Gated RNNs

To remedy the shortcoming of short memory nature of the general RNNs introduced in the last section, *Gated RNNs* was invented in Cho et al. (2014), which included the recently popular long short-term memory (LSTM) networks and those based on the gated recurrent unit (GRU). Gated RNNs use activation function to determine which information is stored in their units for future use; in particular, the trained model can “read” the input temporal sequence of feature vectors, based on which the model decides at some earlier time step t to keep this specific information obtained from the feature vectors for the use at a later time. This stored information can be used later to analyze with the feature vectors recently acquired from the near end of the time series data. For example, if the input text starts with the word “she”, a language processing RNN model could decide to store the information about the gender of the subject so as to interpret properly the possible meaning of “their” to be seen observed in the whole sentence segment.

In a gated RNN, there are different units that make decisions about what informations to store, or when to allow to read, write or erase them. These decisions are learnt from the dataset and are implemented through the concept of “logic” gates. A *gated recurrent unit* (GRU) composes of a reset gate, update gate, cell input and cell output: for instance, in the first layer of a gated recurrent unit (GRU), it takes in two inputs:

1. a feature vector $\mathbf{x}_t =: \mathbf{h}_t^{(0)}$; and
2. a vector from the same or other memory cells in the same (first) layer but from the previous time step, $\mathbf{h}_{t-1}^{(1)}$.

The following Figure 11.4.1 shows a memory cell of a GRU:

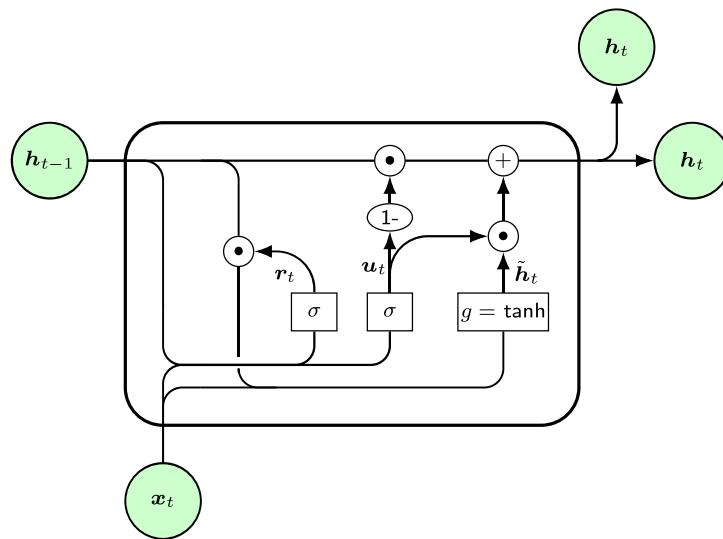


Figure 11.4.1: A GRU memory cell; also see [4].

$$\begin{aligned}
(\text{Reset gate}) \quad & \mathbf{r}_t = \sigma(\mathbf{W}_r \mathbf{x}_t + \mathbf{V}_r \mathbf{h}_{t-1} + \mathbf{b}_r); \\
(\text{Update gate}) \quad & \mathbf{u}_t = \sigma(\mathbf{W}_u \mathbf{x}_t + \mathbf{V}_u \mathbf{h}_{t-1} + \mathbf{b}_u); \\
(\text{Cell input}) \quad & \tilde{\mathbf{h}}_t = \mathbf{g}(\mathbf{W}_{\tilde{h}} \mathbf{x}_t + \mathbf{V}_{\tilde{h}} (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_{\tilde{h}}); \\
(\text{Cell output}) \quad & \mathbf{h}_t = \mathbf{u}_t \odot \tilde{\mathbf{h}}_t + (1 - \mathbf{u}_t) \odot \mathbf{h}_{t-1}.
\end{aligned}
\tag{11.4.1a-d}$$

where in the following, we describe more in detail about the very nature:

1. \mathbf{x}_t stands for the input feature vector at the timestep t ;
2. Matrices $\mathbf{W}_r, \mathbf{W}_u, \mathbf{W}_{\tilde{h}}$ correspond to the weights of inputs \mathbf{x}_t used in the reset gate, update gate, and feedback cell input, respectively;
3. Matrices $\mathbf{V}_r, \mathbf{V}_u, \mathbf{V}_{\tilde{h}}$ correspond to the weights of the connections between the cell output activated value \mathbf{h}_{t-1} being fed back, but one-step delay with the reset gate, update gate, and cell input, respectively;
4. $\mathbf{b}_r, \mathbf{b}_u, \mathbf{b}_{\tilde{h}}$ stands for bias vectors for the reset gate, update gate, and cell input, respectively;
5. Activation functions σ and \mathbf{g} are used for the gates and cell input, respectively. Typically, the logistic function is adopted for σ and the hyperbolic tangent \tanh is used for \mathbf{g} ;
6. **Update Gate:** if $\mathbf{u}_t \approx \mathbf{0}$, then this memory cell, *Memory Gated Unit* (MGU), keeps the value of the cell output \mathbf{h}_{t-1} from the previous time step $t-1$; indeed $\mathbf{h}_t \approx \mathbf{h}_{t-1}$ from (11.4.1d), i.e. \mathbf{h}_{t-1} . It helps to capture the long-term dependence in the time series data;
7. **Reset Gate:** if $\mathbf{u}_t \approx \mathbf{1}$, the value of the MGU is overwritten by a new value $\tilde{\mathbf{h}}_t$, and this is regarded as capturing short-term dependence in time series.

Note that the two gates are mimicking the effect given by a logic gate in electric circuit, namely whenever a threshold is surpassed, an operational signal is fired; otherwise, the state is off. This is also the name of this gated recurrent unit coming from. In particular, in the cell input, we have:

1. if $\mathbf{r}_t = \mathbf{0}$, essentially only information of \mathbf{x}_t is used to calculate $\tilde{\mathbf{h}}_t$;
2. if $\mathbf{r}_t = \mathbf{1}$, \mathbf{h}_{t-1} would interfere, and even override the effects of \mathbf{x}_t especially when $\|\mathbf{V}_{\tilde{h}}\|_2$ is large, that means the previously older memory will reset the state all again and then cover x_t .

A GRU takes an input \mathbf{x}_t and stores it for some time, by then all $\mathbf{u}_t = \mathbf{0}$, for $\tau = 1, \dots, t$, which is somehow equivalent to applying the identity function ($f(x) = x$) to the input. On the other hand, as the gradient of the identity function is a constant vector of “1”, and so it is uniformly bounded away from zero, therefore, when a network with GRUs is trained with backpropagation through time, the gradient does not vanish. To see this effect, let us introduce the backpropagation for this GRU.

11.4.1 Backpropagation through Time

For simplicity, we only consider a 1-layer GRU as before, where we adopt the RNN final output in (11.2.1b) as the output layer, denote the error rates through time be

$$\delta_{t,\tau} := \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_\tau}, \quad \tau = 1, \dots, t,$$

where $\delta_{t,t}$ is given in (11.2.5) as a RNN final output is set. Denote the respective weighted sums \mathbf{z}_t for different gates and inputs be:

$$\begin{aligned}\mathbf{z}_t^r &= \mathbf{W}_r \mathbf{x}_t + \mathbf{V}_r \mathbf{h}_{t-1} + \mathbf{b}_r, & \text{where } \mathbf{r}_t &= \sigma(\mathbf{z}_t^r); \\ \mathbf{z}_t^u &= \mathbf{W}_u \mathbf{x}_t + \mathbf{V}_u \mathbf{h}_{t-1} + \mathbf{b}_u, & \text{where } \mathbf{u}_t &= \sigma(\mathbf{z}_t^u); \\ \mathbf{z}_t^{\tilde{h}} &= \mathbf{W}_{\tilde{h}} \mathbf{x}_t + \mathbf{V}_{\tilde{h}} (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_{\tilde{h}}, & \text{where } \tilde{\mathbf{h}}_t &= \mathbf{g}(\mathbf{z}_t^{\tilde{h}}).\end{aligned}$$

In general, to backpropagate the error rate from $\tau+1$ to τ , we have by simple chain rule:

$$\delta_{t,\tau} := \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_\tau} = \frac{\partial \mathbf{h}_{\tau+1}}{\partial \mathbf{h}_\tau} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_{\tau+1}} = \frac{\partial \mathbf{h}_{\tau+1}}{\partial \mathbf{h}_\tau} \delta_{t,\tau+1}.$$

Recall that the cell output equation in (11.4.1d), we further define a differentiable function $\mathbf{f} : (\mathbf{x}, \mathbf{y}, \mathbf{z}) \mapsto \mathbf{x} \odot \mathbf{y} + (1 - \mathbf{x}) \odot \mathbf{z}$, and also note that both $\mathbf{u}_{\tau+1}$, and $\tilde{\mathbf{h}}_{\tau+1}$ are a function of \mathbf{h}_τ , a simple application of the law of total derivative gives:

$$\begin{aligned}\frac{\partial \mathbf{h}_{\tau+1}}{\partial \mathbf{h}_\tau} &= \frac{\partial \mathbf{f}(\mathbf{u}_{\tau+1}, \tilde{\mathbf{h}}_{\tau+1}, \mathbf{h}_{\tau+1})}{\partial \mathbf{h}_\tau} + \frac{\partial \mathbf{u}_{\tau+1}}{\partial \mathbf{h}_\tau} \frac{\partial \mathbf{f}(\mathbf{u}_{\tau+1}, \tilde{\mathbf{h}}_{\tau+1}, \mathbf{h}_{\tau+1})}{\partial \mathbf{u}_{\tau+1}} + \frac{\partial \tilde{\mathbf{h}}_{\tau+1}}{\partial \mathbf{h}_\tau} \frac{\partial \mathbf{f}(\mathbf{u}_{\tau+1}, \tilde{\mathbf{h}}_{\tau+1}, \mathbf{h}_{\tau+1})}{\partial \tilde{\mathbf{h}}_{\tau+1}} \\ &= \text{diag}\{\mathbf{1} - \mathbf{u}_{\tau+1}\} + \frac{\partial \mathbf{u}_{\tau+1}}{\partial \mathbf{h}_\tau} \frac{\partial \mathbf{h}_{\tau+1}}{\partial \mathbf{u}_{\tau+1}} + \frac{\partial \mathbf{z}_{\tau+1}^{\tilde{h}}}{\partial \mathbf{h}_\tau} \frac{\partial \tilde{\mathbf{h}}_{\tau+1}}{\partial \mathbf{z}_{\tau+1}^{\tilde{h}}} \frac{\partial \mathbf{h}_{\tau+1}}{\partial \tilde{\mathbf{h}}_{\tau+1}},\end{aligned}\tag{11.4.2}$$

where we use (11.1.3) that $(\mathbf{1} - \mathbf{u}_{\tau+1}) \odot \mathbf{h}_\tau = \text{diag}\{\mathbf{1} - \mathbf{u}_{\tau+1}\} \mathbf{h}_\tau$, and then Property 9 in Section 3.4. We further notice that $\mathbf{r}_{\tau+1}$ is also a function of \mathbf{h}_τ in $\mathbf{z}_{\tau+1}^{\tilde{h}}$, therefore in particular, the law of total derivative give:

$$\frac{\partial \mathbf{z}_{\tau+1}^{\tilde{h}}}{\partial \mathbf{h}_\tau} = \frac{\partial (\mathbf{V}_{\tilde{h}} \text{diag}\{\mathbf{h}_\tau\} \mathbf{r}_{\tau+1})}{\partial \mathbf{h}_\tau} = \frac{\partial (\mathbf{V}_{\tilde{h}} \text{diag}\{\mathbf{r}_{\tau+1}\} \mathbf{h}_\tau)}{\partial \mathbf{h}_\tau} = \frac{\partial \mathbf{r}_{\tau+1}}{\partial \mathbf{h}_\tau} (\mathbf{V}_{\tilde{h}} \text{diag}\{\mathbf{h}_\tau\})^\top + (\mathbf{V}_{\tilde{h}} \text{diag}\{\mathbf{r}_{\tau+1}\})^\top.$$

Therefore, (11.4.2) can be further computed as:

$$\begin{aligned}\frac{\partial \mathbf{h}_{\tau+1}}{\partial \mathbf{h}_\tau} &= \text{diag}\{\mathbf{1} - \mathbf{u}_{\tau+1}\} + \mathbf{V}_u^\top \text{diag}\{\sigma'(\mathbf{z}_{\tau+1}^u)\} \text{diag}\{\tilde{\mathbf{h}}_{\tau+1} - \mathbf{h}_{\tau+1}\} \\ &\quad + \left[\frac{\partial \mathbf{r}_{\tau+1}}{\partial \mathbf{h}_\tau} (\mathbf{V}_{\tilde{h}} \text{diag}\{\mathbf{h}_\tau\})^\top + (\mathbf{V}_{\tilde{h}} \text{diag}\{\mathbf{r}_{\tau+1}\})^\top \right] \text{diag}\{\mathbf{g}'(\mathbf{z}_{\tau+1}^{\tilde{h}})\} \text{diag}\{\mathbf{u}_{\tau+1}\} \\ &= \text{diag}\{\mathbf{1} - \mathbf{u}_{\tau+1}\} + \mathbf{V}_u^\top \text{diag}\left\{\sigma'(\mathbf{z}_{\tau+1}^u) \odot (\tilde{\mathbf{h}}_{\tau+1} - \mathbf{h}_{\tau+1})\right\} \\ &\quad + \left[\mathbf{V}_r^\top \text{diag}\left\{\sigma'(\mathbf{z}_{\tau+1}^r) \odot \mathbf{h}_\tau\right\} + \text{diag}\{\mathbf{r}_{\tau+1}\} \right] \mathbf{V}_{\tilde{h}}^\top \text{diag}\left\{\mathbf{g}'(\mathbf{z}_{\tau+1}^{\tilde{h}}) \odot \mathbf{u}_{\tau+1}\right\}.\end{aligned}$$

Therefore, for $\tau = 1, \dots, t-1$, we have

$$\begin{aligned}\delta_{t,\tau} &= \delta_{t,\tau+1} \odot (\mathbf{1} - \mathbf{u}_{\tau+1}) + \mathbf{V}_u^\top \left(\delta_{t,\tau+1} \odot \sigma'(\mathbf{z}_{\tau+1}^u) \odot (\tilde{\mathbf{h}}_{\tau+1} - \mathbf{h}_{\tau+1}) \right) \\ &\quad + \left[\mathbf{V}_r^\top \text{diag}\left\{\sigma'(\mathbf{z}_{\tau+1}^r) \odot \mathbf{h}_\tau\right\} + \text{diag}\{\mathbf{r}_{\tau+1}\} \right] \mathbf{V}_{\tilde{h}}^\top \left(\delta_{t,\tau+1} \odot \mathbf{g}'(\mathbf{z}_{\tau+1}^{\tilde{h}}) \odot \mathbf{u}_{\tau+1} \right).\end{aligned}\tag{11.4.3}$$

Consider the parameter updates in the cell input,

$$\begin{aligned}\frac{\partial \mathcal{E}}{\partial \mathbf{W}_{\tilde{h}}} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{W}_{\tilde{h}}} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{h}_\tau}{\partial \mathbf{W}_{\tilde{h}}} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \tilde{\mathbf{h}}_\tau}{\partial \mathbf{W}_{\tilde{h}}} \frac{\partial \mathbf{h}_\tau}{\partial \tilde{\mathbf{h}}_\tau} \delta_{t,\tau} \\ &= \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial (\mathbf{W}_{\tilde{h}} \mathbf{x}_\tau)}{\partial \mathbf{W}_{\tilde{h}}} \text{diag}\{\mathbf{g}'(\mathbf{z}_\tau^{\tilde{h}})\} \text{diag}(\mathbf{u}_\tau) \delta_{t,\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \{\delta_{t,\tau} \odot \mathbf{g}'(\mathbf{z}_\tau^{\tilde{h}}) \odot \mathbf{u}_\tau\} \mathbf{x}_\tau^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{V}_{\tilde{h}}} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{V}_{\tilde{h}}} = \sum_{t=1}^D \sum_{\tau=1}^t \{\delta_{t,\tau} \odot \mathbf{g}'(\mathbf{z}_\tau^{\tilde{h}}) \odot \mathbf{u}_\tau\} \mathbf{h}_{\tau-1}^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{b}_{\tilde{h}}} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{b}_{\tilde{h}}} = \sum_{t=1}^D \sum_{\tau=1}^t \{\delta_{t,\tau} \odot \mathbf{g}'(\mathbf{z}_\tau^{\tilde{h}}) \odot \mathbf{u}_\tau\}.\end{aligned}$$

While for the parameter updates in the update gate, likewise we have:

$$\begin{aligned}\frac{\partial \mathcal{E}}{\partial \mathbf{W}_u} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{W}_u} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{h}_\tau}{\partial \mathbf{W}_u} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{u}_\tau}{\partial \mathbf{W}_u} \frac{\partial \mathbf{h}_\tau}{\partial \mathbf{u}_\tau} \delta_{t,\tau} \\ &= \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial (\mathbf{W}_u \mathbf{x}_\tau)}{\partial \mathbf{W}_u} \text{diag}\{\boldsymbol{\sigma}'(\mathbf{z}_\tau^u)\} \text{diag}\{\tilde{\mathbf{h}}_\tau - \mathbf{h}_{\tau-1}\} \delta_{t,\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \delta_{t,\tau} \odot \boldsymbol{\sigma}'(\mathbf{z}_\tau^u) \odot (\tilde{\mathbf{h}}_\tau - \mathbf{h}_{\tau-1}) \right\} \mathbf{x}_\tau^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{V}_u} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{V}_u} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \delta_{t,\tau} \odot \boldsymbol{\sigma}'(\mathbf{z}_\tau^u) \odot (\tilde{\mathbf{h}}_\tau - \mathbf{h}_{\tau-1}) \right\} \mathbf{h}_{\tau-1}^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{b}_u} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{b}_u} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \delta_{t,\tau} \odot \boldsymbol{\sigma}'(\mathbf{z}_\tau^u) \odot (\tilde{\mathbf{h}}_\tau - \mathbf{h}_{\tau-1}) \right\}.\end{aligned}$$

Finally, for the parameter updates in the reset gate,

$$\begin{aligned}\frac{\partial \mathcal{E}}{\partial \mathbf{W}_r} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{W}_r} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{h}_\tau}{\partial \mathbf{W}_r} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{r}_\tau}{\partial \mathbf{W}_r} \frac{\partial \tilde{\mathbf{h}}_\tau}{\partial \mathbf{r}_\tau} \frac{\partial \mathbf{h}_\tau}{\partial \tilde{\mathbf{h}}_\tau} \delta_{t,\tau} \\ &= \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial (\mathbf{W}_r \mathbf{x}_\tau)}{\partial \mathbf{W}_r} \text{diag}\{\boldsymbol{\sigma}'(\mathbf{z}_\tau^r)\} \mathbf{V}_{\tilde{h}}^\top \text{diag}\{\mathbf{g}'(\mathbf{z}_\tau^{\tilde{h}})\} \text{diag}\{\mathbf{u}_\tau\} \delta_{t,\tau} \\ &= \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \boldsymbol{\sigma}'(\mathbf{z}_\tau^r) \odot \left(\mathbf{V}_{\tilde{h}}^\top \{\delta_{t,\tau} \odot \mathbf{g}'(\mathbf{z}_\tau^{\tilde{h}}) \odot \mathbf{u}_\tau\} \right) \right\} \mathbf{x}_\tau^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{V}_r} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{V}_r} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \boldsymbol{\sigma}'(\mathbf{z}_\tau^r) \odot \left(\mathbf{V}_{\tilde{h}}^\top \{\delta_{t,\tau} \odot \mathbf{g}'(\mathbf{z}_\tau^{\tilde{h}}) \odot \mathbf{u}_\tau\} \right) \right\} \mathbf{h}_{\tau-1}^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{b}_r} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{b}_r} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \boldsymbol{\sigma}'(\mathbf{z}_\tau^r) \odot \left(\mathbf{V}_{\tilde{h}}^\top \{\delta_{t,\tau} \odot \mathbf{g}'(\mathbf{z}_\tau^{\tilde{h}}) \odot \mathbf{u}_\tau\} \right) \right\}.\end{aligned}$$

From (11.4.3), we observe that the cell error rate can be expressed in the form:

$$\delta_{t,\tau} = (\mathbf{1} - \mathbf{u}_{\tau+1}) \odot \delta_{t,\tau+1} + \mathbf{A}_{\tau+1}(\mathbf{V}_u) + \mathbf{B}_{\tau+1}(\mathbf{V}_{\tilde{h}}, \mathbf{V}_r),$$

where the first term does not depend on the weights \mathbf{V} 's, and

$$\begin{aligned}\mathbf{A}_{\tau+1}(\mathbf{V}_u) &= \mathbf{V}_u^\top \text{diag}\left\{ \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^u) \odot (\tilde{\mathbf{h}}_{\tau+1} - \mathbf{h}_{\tau+1}) \right\}; \\ \mathbf{B}_{\tau+1}(\mathbf{V}_{\tilde{h}}, \mathbf{V}_r) &= \left[\mathbf{V}_r^\top \text{diag}\left\{ \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^r) \odot \mathbf{h}_\tau \right\} + \text{diag}\{\mathbf{r}_{\tau+1}\} \right] \mathbf{V}_{\tilde{h}}^\top \left(\delta_{t,\tau+1} \odot \mathbf{g}'(\mathbf{z}_{\tau+1}^{\tilde{h}}) \odot \mathbf{u}_{\tau+1} \right).\end{aligned}$$

Therefore, if the value of the update gate \mathbf{u}_τ is close to $\mathbf{0}$ for all $\tau = t-1, \dots, 1$, then with $\boldsymbol{\sigma}$ being the logistic activation function, we have $\boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^u) \approx \mathbf{0}$, and hence the cell error rate in (11.4.3) is simply $\delta_{t,\tau} \approx \delta_{t,\tau+1}$. In general, the error rates of a fully-trained network are close to zero such that the updates of the parameters in the network are small. However, when suffering from the vanishing gradient problem, the error rates diminish to zero in the process of backpropagating through time even when the model is not fully-trained, *i.e.* the optimal parameters are not yet attained. With $\mathbf{u}_\tau \approx \mathbf{0}$ in the gated RNN, we have a stable error rates at all timestep τ as $\delta_{t,\tau} \approx \delta_{t,\tau+1}$ across $\tau = 1, \dots, t-1$, and hence, as long as $\delta_{t,t}$ is significantly different from zero, the gated RNN avoids the vanishing gradient problem.

11.4.2 Example of Gated RNN: LSTM

Long Short-Term Memory (LSTM), first proposed in Hochreiter and Schmidhuber (1997), makes use of memory cells (central processing and storage units) that allow a constant propagation flow of error rates

through layers and time during the training of the network; in this way, as GRU, LSTM can avoid the vanishing gradient problem, and hence it allows for a *uniform credit assignment*, which means the backpropagation of error rates to inputs without scaling to be described further in Subsection 11.4.3. The following Figure 11.4.2 shows a memory cell of LSTM.

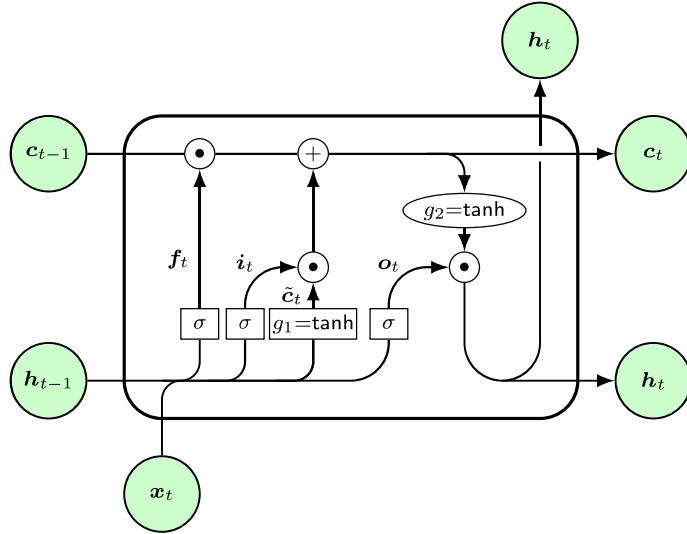


Figure 11.4.2: LSTM memory cell; also see [4].

(Forget gate)

$$f_t = \sigma(W_f x_t + V_f h_{t-1} + b_f); \quad (11.4.4a)$$

(Input gate)

$$i_t = \sigma(W_i x_t + V_i h_{t-1} + b_i); \quad (11.4.4b)$$

(Cell input)

$$\tilde{c}_t = g_1(W_{\tilde{c}} x_t + V_{\tilde{c}} h_{t-1} + b_{\tilde{c}}); \quad (11.4.4c)$$

(Output gate)

$$o_t = \sigma(W_o x_t + V_o h_{t-1} + b_o); \quad (11.4.4d)$$

(Cell state)

$$c_t = i_t \odot \tilde{c}_t + f_t \odot c_{t-1}; \quad (11.4.4e)$$

(Cell output)

$$h_t = o_t \odot g_2(c_t); \quad (11.4.4f)$$

where

1. for this memory cell, x_t stands for the input feature vector at the present timestep t ;
2. matrices $W_f, W_i, W_{\tilde{c}}, W_o$ correspond to the weights of inputs x_t used for the forget gate, input gate, cell input and output gate, respectively;
3. matrices $V_f, V_i, V_{\tilde{c}}, V_o$ correspond to the weights of the cell output activated value h_{t-1} , fed from the one-step delay of itself for the forget gate, input gate, cell input and output gate, respectively;
4. vectors $b_f, b_i, b_{\tilde{c}}, b_o$ correspond to the biases of the forget gate, input gate, cell input, and output gate, respectively;
5. activation functions g_1, g_2, σ are used for the cell input, cell state and all other gates, respectively. Typically, hyperbolic tangent \tanh is chosen for g_1 and g_2 ; while logistic function can be adopted for σ ;
6. **Forget Gate f_t :** decides what the last timestep cell input information should be dropped or kept. As an image value under a sigmoid type activation function, ranging in $(0, 1)$, its value closer to 0 means to forget, while if it is closer to 1, it means to keep;

7. **Input Gate i_t :** decides whether the input should be updated. As a sigmoid type activation function taking values in $(0, 1)$, if its value is closer to 0 that means no update required but just using the last timestep; otherwise if it is closer to 1, one has to update by using \tilde{c}_t ; and
8. **Output Gate o_t :** decides whether the current cell state will previously stored input value provide a non-zero output of the cell, which serves as an input for the next hidden layer.
9. with the input gate i and forget gate f taking values in $(0, 1)$, the cell state c_t represents the long term memory. If the forget gate is close to 1, then we keep all the cell state “memories” and further accumulate it with the (short-term information) cell input \tilde{c} through the input gate when forward propagating through time.
10. with g_2 taking values in $(-1, 1)$, the cell output h_t represents the short-term memory.

11.4.2.1 Backpropagation through Time in LSTM

For illustration of the main idea, we only consider a 1-layer LSTM, where we adopt the RNN final output in (11.2.1b) as the output layer. As there are two products generated by the memory cell, namely the cell output and the cell state, we denote the two corresponding error rates be

$$\delta_{t,\tau} := \frac{\partial \mathcal{E}_t}{\partial h_\tau} \quad \text{and} \quad \lambda_{t,\tau} := \frac{\partial \mathcal{E}_t}{\partial c_\tau}, \quad \tau = 1, \dots, t,$$

where $\delta_{t,t}$ is given in (11.2.5) as a RNN final output is set. Consider the cell output equation:

$$h_t = o_t \odot g_2(c_t).$$

the respective partial derivatives of the error with respect to o_t and c_t are:

$$\frac{\partial \mathcal{E}_t}{\partial o_t} = \frac{\partial h_t}{\partial o_t} \frac{\partial \mathcal{E}_t}{\partial h_t} = \text{diag}\{g_2(c_t)\} \delta_{t,t} = \delta_{t,t} \odot g_2(c_t); \quad (11.4.5)$$

$$\frac{\partial \mathcal{E}_t}{\partial c_t} = \frac{\partial h_t}{\partial c_t} \frac{\partial \mathcal{E}_t}{\partial h_t} = \text{diag}\{o_t \odot g'_2(c_t)\} \delta_{t,t} = \delta_{t,t} \odot o_t \odot g'_2(c_t) =: \lambda_{t,t}. \quad (11.4.6)$$

Denote the respective weighted sums z_t for different gates and inputs by:

$$z_t^f = W_f x_t + V_f h_{t-1} + b_f, \quad \text{where } f_t = \sigma(z_t^f);$$

$$z_t^i = W_i x_t + V_i h_{t-1} + b_i, \quad \text{where } i_t = \sigma(z_t^i);$$

$$z_t^{\tilde{c}} = W_{\tilde{c}} x_t + V_{\tilde{c}} h_{t-1} + b_{\tilde{c}}, \quad \text{where } \tilde{c}_t = g_1(z_t^{\tilde{c}});$$

$$z_t^o = W_o x_t + V_o h_{t-1} + b_o, \quad \text{where } o_t = \sigma(z_t^o).$$

To backpropagate the cell error rate $\lambda_{t,\tau+1}$ from $\tau+1$ to τ , for $\tau = t-1, \dots, 1$,

$$\lambda_{t,\tau} = \frac{\partial \mathcal{E}_t}{\partial c_\tau} = \frac{\partial c_{\tau+1}}{\partial c_\tau} \frac{\partial \mathcal{E}_t}{\partial c_{\tau+1}} = \frac{\partial c_{\tau+1}}{\partial c_\tau} \lambda_{t,\tau+1}. \quad (11.4.7)$$

Recall (11.4.4e) at $\tau+1$, we have $c_{\tau+1} = i_{\tau+1} \odot \tilde{c}_{\tau+1} + f_{\tau+1} \odot c_\tau$. On the other hand, at the timestep τ in (11.4.4f), we have $h_\tau = o_\tau \odot g_2(c_\tau)$ with which all $i_{\tau+1}$, $\tilde{c}_{\tau+1}$, and $f_{\tau+1}$ are functions of c_τ . Therefore, by the law of total derivative, similar to deriving (11.4.2),

$$\begin{aligned} \frac{\partial c_{\tau+1}}{\partial c_\tau} &= \frac{\partial h_\tau}{\partial c_\tau} \frac{\partial i_{\tau+1}}{\partial h_\tau} \frac{\partial c_{\tau+1}}{\partial i_{\tau+1}} + \frac{\partial h_\tau}{\partial c_\tau} \frac{\partial \tilde{c}_{\tau+1}}{\partial h_\tau} \frac{\partial c_{\tau+1}}{\partial \tilde{c}_{\tau+1}} + \frac{\partial h_\tau}{\partial c_\tau} \frac{\partial f_{\tau+1}}{\partial h_\tau} \frac{\partial c_{\tau+1}}{\partial f_{\tau+1}} + \text{diag}\{f_{\tau+1}\} \\ &= \text{diag}\{o_\tau \odot g'_2(c_\tau)\} \left\{ \frac{\partial i_{\tau+1}}{\partial h_\tau} \text{diag}\{\tilde{c}_{\tau+1}\} + \frac{\partial \tilde{c}_{\tau+1}}{\partial h_\tau} \text{diag}\{i_{\tau+1}\} + \frac{\partial f_{\tau+1}}{\partial h_\tau} \text{diag}\{c_\tau\} \right\} + \text{diag}\{f_{\tau+1}\} \\ &= \text{diag}\{o_\tau \odot g'_2(c_\tau)\} \left\{ V_i^\top \text{diag}\{\sigma'(z_{\tau+1}^i) \odot \tilde{c}_{\tau+1}\} + V_{\tilde{c}}^\top \text{diag}\{g'_1(z_{\tau+1}^{\tilde{c}}) \odot i_{\tau+1}\} + V_f^\top \text{diag}\{\sigma'(z_{\tau+1}^f) \odot c_\tau\} \right\} \\ &\quad + \text{diag}\{f_{\tau+1}\}. \end{aligned}$$

Therefore, for $\tau = 1, \dots, t-1$, the cell error rate in (11.4.7) can be further computed as:

$$\begin{aligned}\boldsymbol{\lambda}_{t,\tau} &= \text{diag} \left\{ \mathbf{o}_\tau \odot \mathbf{g}'_2(\mathbf{c}_\tau) \right\} \left[\mathbf{V}_i^\top \left\{ \boldsymbol{\lambda}_{t,\tau+1} \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^i) \odot \tilde{\mathbf{c}}_{\tau+1} \right\} + \mathbf{V}_{\tilde{c}}^\top \left\{ \boldsymbol{\lambda}_{t,\tau+1} \odot \mathbf{g}'_1(\mathbf{z}_{\tau+1}^{\tilde{c}}) \odot \mathbf{i}_{\tau+1} \right\} \right. \\ &\quad \left. + \mathbf{V}_f^\top \left\{ \boldsymbol{\lambda}_{t,\tau+1} \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^f) \odot \mathbf{c}_\tau \right\} \right] + \mathbf{f}_{\tau+1} \odot \boldsymbol{\lambda}_{t,\tau+1}.\end{aligned}\quad (11.4.8)$$

Consider the parameter updates in the forget gate,

$$\begin{aligned}\frac{\partial \mathcal{E}}{\partial \mathbf{W}_f} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{W}_f} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{c}_\tau}{\partial \mathbf{W}_f} \frac{\partial \mathcal{E}_t}{\partial \mathbf{c}_\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{f}_\tau}{\partial \mathbf{W}_f} \frac{\partial \mathbf{c}_\tau}{\partial \mathbf{f}_\tau} \boldsymbol{\lambda}_{t,\tau} \\ &= \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial (\mathbf{W}_f \mathbf{x}_\tau)}{\partial \mathbf{W}_f} \text{diag} \left\{ \boldsymbol{\sigma}'(\mathbf{z}_\tau^f) \right\} \left(\text{diag} \left\{ \mathbf{c}_{\tau-1} \right\} \boldsymbol{\lambda}_{t,\tau} \right) = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \boldsymbol{\lambda}_{t,\tau} \odot \mathbf{c}_{\tau-1} \odot \boldsymbol{\sigma}'(\mathbf{z}_\tau^f) \right\} \mathbf{x}_\tau^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{V}_f} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{V}_f} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \boldsymbol{\lambda}_{t,\tau} \odot \mathbf{c}_{\tau-1} \odot \boldsymbol{\sigma}'(\mathbf{z}_\tau^f) \right\} \mathbf{h}_{\tau-1}^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{b}_f} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{b}_f} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \boldsymbol{\lambda}_{t,\tau} \odot \mathbf{c}_{\tau-1} \odot \boldsymbol{\sigma}'(\mathbf{z}_\tau^f) \right\}.\end{aligned}$$

Likewise, for the parameter updates in the input gate,

$$\begin{aligned}\frac{\partial \mathcal{E}}{\partial \mathbf{W}_i} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{W}_i} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{c}_\tau}{\partial \mathbf{W}_i} \frac{\partial \mathcal{E}_t}{\partial \mathbf{c}_\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{i}_\tau}{\partial \mathbf{W}_i} \frac{\partial \mathbf{c}_\tau}{\partial \mathbf{i}_\tau} \boldsymbol{\lambda}_{t,\tau} \\ &= \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial (\mathbf{W}_i \mathbf{x}_\tau)}{\partial \mathbf{W}_i} \text{diag} \left\{ \boldsymbol{\sigma}'(\mathbf{z}_\tau^i) \right\} \text{diag} \left\{ \tilde{\mathbf{c}}_\tau \right\} \boldsymbol{\lambda}_{t,\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \boldsymbol{\lambda}_{t,\tau} \odot \tilde{\mathbf{c}}_\tau \odot \boldsymbol{\sigma}'(\mathbf{z}_\tau^i) \right\} \mathbf{x}_\tau^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{V}_i} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{V}_i} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \boldsymbol{\lambda}_{t,\tau} \odot \tilde{\mathbf{c}}_\tau \odot \boldsymbol{\sigma}'(\mathbf{z}_\tau^i) \right\} \mathbf{h}_{\tau-1}^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{b}_i} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{b}_i} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \boldsymbol{\lambda}_{t,\tau} \odot \tilde{\mathbf{c}}_\tau \odot \boldsymbol{\sigma}'(\mathbf{z}_\tau^i) \right\}.\end{aligned}$$

Similarly, for the parameter updates in the cell input,

$$\begin{aligned}\frac{\partial \mathcal{E}}{\partial \mathbf{W}_{\tilde{c}}} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{W}_{\tilde{c}}} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{c}_\tau}{\partial \mathbf{W}_{\tilde{c}}} \frac{\partial \mathcal{E}_t}{\partial \mathbf{c}_\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \tilde{\mathbf{c}}_\tau}{\partial \mathbf{W}_{\tilde{c}}} \frac{\partial \mathbf{c}_\tau}{\partial \tilde{\mathbf{c}}_\tau} \boldsymbol{\lambda}_{t,\tau} \\ &= \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial (\mathbf{W}_{\tilde{c}} \mathbf{x}_\tau)}{\partial \mathbf{W}_{\tilde{c}}} \text{diag} \left\{ \mathbf{g}'_1(\mathbf{z}_\tau^{\tilde{c}}) \right\} \text{diag} \left\{ \mathbf{i}_\tau \right\} \boldsymbol{\lambda}_{t,\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \boldsymbol{\lambda}_{t,\tau} \odot \mathbf{i}_\tau \odot \mathbf{g}'_1(\mathbf{z}_\tau^{\tilde{c}}) \right\} \mathbf{x}_\tau^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{V}_{\tilde{c}}} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{V}_{\tilde{c}}} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \boldsymbol{\lambda}_{t,\tau} \odot \mathbf{i}_\tau \odot \mathbf{g}'_1(\mathbf{z}_\tau^{\tilde{c}}) \right\} \mathbf{h}_{\tau-1}^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{b}_{\tilde{c}}} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{b}_{\tilde{c}}} = \sum_{t=1}^D \sum_{\tau=1}^t \left\{ \boldsymbol{\lambda}_{t,\tau} \odot \mathbf{i}_\tau \odot \mathbf{g}'_1(\mathbf{z}_\tau^{\tilde{c}}) \right\}.\end{aligned}$$

Finally, there is still an output gate to be dealt, which involves more argument as follows. The next step is to backpropagate the state error rate $\delta_{t,\tau+1}$ from $\tau+1$ to τ , i.e.

$$\delta_{t,\tau} = \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_\tau} = \frac{\partial \mathbf{h}_{\tau+1}}{\partial \mathbf{h}_\tau} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_{\tau+1}} = \frac{\partial \mathbf{h}_{\tau+1}}{\partial \mathbf{h}_\tau} \delta_{t,\tau+1}.$$

Note that $\mathbf{o}_{\tau+1}$ and $\mathbf{c}_{\tau+1}$ are functions of \mathbf{h}_τ from (11.4.4d) and (11.4.4e), which is further dealt to (11.4.4a), (11.4.4b) and (11.4.4c). From (11.4.4f) at $\tau+1$ that $\mathbf{h}_{\tau+1} = \mathbf{o}_{\tau+1} \odot \mathbf{g}_2(\mathbf{c}_{\tau+1})$, by law of total derivative, we

have

$$\begin{aligned}\frac{\partial \mathbf{h}_{\tau+1}}{\partial \mathbf{h}_\tau} &= \frac{\partial \mathbf{o}_{\tau+1}}{\partial \mathbf{h}_\tau} \frac{\partial \mathbf{h}_{\tau+1}}{\partial \mathbf{o}_{\tau+1}} + \frac{\partial \mathbf{c}_{\tau+1}}{\partial \mathbf{h}_\tau} \frac{\partial \mathbf{h}_{\tau+1}}{\partial \mathbf{c}_{\tau+1}} = \mathbf{V}_o^\top \text{diag}\{\boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^o)\} \text{diag}\{\mathbf{1} \odot \mathbf{g}_2(\mathbf{c}_{\tau+1})\} + \frac{\partial \mathbf{c}_{\tau+1}}{\partial \mathbf{h}_\tau} \text{diag}\{\mathbf{o}_{\tau+1} \odot \mathbf{g}'_2(\mathbf{c}_{\tau+1})\} \\ &= \mathbf{V}_o^\top \text{diag}\{\mathbf{g}_2(\mathbf{c}_{\tau+1}) \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^o)\} + \frac{\partial \mathbf{c}_{\tau+1}}{\partial \mathbf{h}_\tau} \text{diag}\{\mathbf{o}_{\tau+1} \odot \mathbf{g}'_2(\mathbf{c}_{\tau+1})\}.\end{aligned}$$

To compute $\frac{\partial \mathbf{c}_{\tau+1}}{\partial \mathbf{h}_\tau}$, we first recall that $\mathbf{f}_{\tau+1}$, $\mathbf{i}_{\tau+1}$, and $\tilde{\mathbf{c}}_{\tau+1}$ are functions of \mathbf{h}_τ from (11.4.4a), (11.4.4b), and (11.4.4c) respectively. Then, applying the law of total derivative to (11.4.4e), we obtain

$$\begin{aligned}\frac{\partial \mathbf{c}_{\tau+1}}{\partial \mathbf{h}_\tau} &= \frac{\partial \mathbf{i}_{\tau+1}}{\partial \mathbf{h}_\tau} \frac{\partial \mathbf{c}_{\tau+1}}{\partial \mathbf{i}_{\tau+1}} + \frac{\partial \mathbf{f}_{\tau+1}}{\partial \mathbf{h}_\tau} \frac{\partial \mathbf{c}_{\tau+1}}{\partial \mathbf{f}_{\tau+1}} + \frac{\partial \tilde{\mathbf{c}}_{\tau+1}}{\partial \mathbf{h}_\tau} \frac{\partial \mathbf{c}_{\tau+1}}{\partial \tilde{\mathbf{c}}_{\tau+1}} \\ &= \mathbf{V}_i^\top \text{diag}\{\boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^i)\} \text{diag}\{\tilde{\mathbf{c}}_{\tau+1}\} + \mathbf{V}_f^\top \text{diag}\{\boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^f)\} \text{diag}\{\mathbf{c}_\tau\} + \mathbf{V}_{\tilde{c}}^\top \text{diag}\{\boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^{\tilde{c}})\} \text{diag}\{\mathbf{i}_{\tau+1}\} \\ &= \mathbf{V}_i^\top \text{diag}\{\tilde{\mathbf{c}}_{\tau+1} \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^i)\} + \mathbf{V}_f^\top \text{diag}\{\mathbf{c}_\tau \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^f)\} + \mathbf{V}_{\tilde{c}}^\top \text{diag}\{\mathbf{i}_{\tau+1} \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^{\tilde{c}})\}.\end{aligned}$$

Therefore, for $\tau = 1, \dots, t-1$, we have

$$\begin{aligned}\delta_{t,\tau} &= \mathbf{V}_o^\top \left\{ \delta_{t,\tau+1} \odot \mathbf{g}_2(\mathbf{c}_{\tau+1}) \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^o) \right\} + \left[\mathbf{V}_i^\top \text{diag}\{\tilde{\mathbf{c}}_{\tau+1} \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^i)\} \right. \\ &\quad \left. + \mathbf{V}_f^\top \text{diag}\{\mathbf{c}_\tau \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^f)\} + \mathbf{V}_{\tilde{c}}^\top \text{diag}\{\mathbf{i}_{\tau+1} \odot \mathbf{g}'_1(\mathbf{z}_{\tau+1}^{\tilde{c}})\} \right] \left\{ \delta_{t,\tau+1} \odot \mathbf{o}_{\tau+1} \odot \mathbf{g}'_2(\mathbf{c}_{\tau+1}) \right\} \\ &= \mathbf{V}_o^\top \left\{ \delta_{t,\tau+1} \odot \mathbf{g}_2(\mathbf{c}_{\tau+1}) \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^o) \right\} + \mathbf{V}_i^\top \left\{ \delta_{t,\tau+1} \odot \tilde{\mathbf{c}}_{\tau+1} \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^i) \odot \mathbf{o}_{\tau+1} \odot \mathbf{g}'_2(\mathbf{c}_{\tau+1}) \right\} \\ &\quad + \mathbf{V}_f^\top \left\{ \delta_{t,\tau+1} \odot \mathbf{c}_\tau \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^f) \odot \mathbf{o}_{\tau+1} \odot \mathbf{g}'_2(\mathbf{c}_{\tau+1}) \right\} + \mathbf{V}_{\tilde{c}}^\top \left\{ \delta_{t,\tau+1} \odot \mathbf{i}_{\tau+1} \odot \mathbf{g}'_1(\mathbf{z}_{\tau+1}^{\tilde{c}}) \odot \mathbf{o}_{\tau+1} \odot \mathbf{g}'_2(\mathbf{c}_{\tau+1}) \right\} \\ &= \begin{pmatrix} \mathbf{V}_o^\top & \mathbf{V}_i^\top & \mathbf{V}_f^\top & \mathbf{V}_{\tilde{c}}^\top \end{pmatrix} \begin{pmatrix} \delta_{t,\tau+1} \odot \mathbf{g}_2(\mathbf{c}_{\tau+1}) \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^o) \\ \delta_{t,\tau+1} \odot \tilde{\mathbf{c}}_{\tau+1} \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^i) \odot \mathbf{o}_{\tau+1} \odot \mathbf{g}'_2(\mathbf{c}_{\tau+1}) \\ \delta_{t,\tau+1} \odot \mathbf{c}_\tau \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^f) \odot \mathbf{o}_{\tau+1} \odot \mathbf{g}'_2(\mathbf{c}_{\tau+1}) \\ \delta_{t,\tau+1} \odot \mathbf{i}_{\tau+1} \odot \mathbf{g}'_1(\mathbf{z}_{\tau+1}^{\tilde{c}}) \odot \mathbf{o}_{\tau+1} \odot \mathbf{g}'_2(\mathbf{c}_{\tau+1}) \end{pmatrix}.\end{aligned}$$

Eventually, for the parameter updates in the output gate, we see:

$$\begin{aligned}\frac{\partial \mathcal{E}}{\partial \mathbf{W}_o} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{W}_o} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial \mathbf{h}_\tau}{\partial \mathbf{W}_o} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_\tau} = \sum_{t=1}^D \sum_{\tau=1}^t \frac{\partial (\mathbf{W}_o \mathbf{x}_\tau)}{\partial \mathbf{W}_o} \text{diag}\{\boldsymbol{\sigma}'(\mathbf{z}_\tau^o)\} \text{diag}\{\mathbf{g}_2(\mathbf{c}_\tau)\} \delta_{t,\tau} \\ &= \sum_{t=1}^D \sum_{\tau=1}^t \{\delta_{t,\tau} \odot \mathbf{g}_2(\mathbf{c}_\tau) \odot \boldsymbol{\sigma}'(\mathbf{z}_\tau^o)\} \mathbf{x}_\tau^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{V}_o} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{V}_o} = \sum_{t=1}^D \sum_{\tau=1}^t \{\delta_{t,\tau} \odot \mathbf{g}_2(\mathbf{c}_\tau) \odot \boldsymbol{\sigma}'(\mathbf{z}_\tau^o)\} \mathbf{h}_{\tau-1}^\top, \\ \frac{\partial \mathcal{E}}{\partial \mathbf{b}_o} &= \sum_{t=1}^D \frac{\partial \mathcal{E}_t}{\partial \mathbf{b}_o} = \sum_{t=1}^D \sum_{\tau=1}^t \{\delta_{t,\tau} \odot \mathbf{g}_2(\mathbf{c}_\tau) \odot \boldsymbol{\sigma}'(\mathbf{z}_\tau^o)\}.\end{aligned}$$

11.4.3 Uniform Credit Assignment

Uniform credit assignment stands for the property that LSTM can consider all input information at each phase of learning, no matter where they are located in the input temporal sequence. From (11.4.8), we observe that the cell error rate can be written in the form:

$$\lambda_{t,\tau} = \text{diag}\{\mathbf{o}_\tau \odot \mathbf{g}'_2(\mathbf{c}_\tau)\} \left(\mathbf{A}_{\tau+1}(\mathbf{V}_i) + \mathbf{B}_{\tau+1}(\mathbf{V}_f) + \mathbf{C}_{\tau+1}(\mathbf{V}_{\tilde{c}}) \right) + \mathbf{f}_{\tau+1} \lambda_{t,\tau+1},$$

where the last term does not depend on the weights \mathbf{V} , and

$$\begin{aligned}\mathbf{A}_{\tau+1}(\mathbf{V}_i) &= \mathbf{V}_i^\top \left\{ \boldsymbol{\lambda}_{t,\tau+1} \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^i) \odot \tilde{\mathbf{c}}_{\tau+1} \right\}; \\ \mathbf{B}_{\tau+1}(\mathbf{V}_f) &= \mathbf{V}_f^\top \left\{ \boldsymbol{\lambda}_{t,\tau+1} \odot \boldsymbol{\sigma}'(\mathbf{z}_{\tau+1}^f) \odot \mathbf{c}_\tau \right\}; \\ \mathbf{C}_{\tau+1}(\mathbf{V}_c) &= \mathbf{V}_c^\top \left\{ \boldsymbol{\lambda}_{t,\tau+1} \odot \mathbf{g}'_1(\mathbf{z}_{\tau+1}^c) \odot \mathbf{i}_{\tau+1} \right\}.\end{aligned}$$

Therefore, if the value of the forget gate \mathbf{f} is close to $\mathbf{1}$ and that of the input gate \mathbf{i} is close to $\mathbf{0}$, then both \mathbf{z}^f and \mathbf{z}^i take extreme values, very large in positive and in negative, respectively. Therefore, $\boldsymbol{\sigma}'(\mathbf{z}^f) \approx \mathbf{0}$ and $\boldsymbol{\sigma}'(\mathbf{z}^i) \approx \mathbf{0}$, as a result $\mathbf{A}_{\tau+1}(\mathbf{V}_i)$, $\mathbf{B}_{\tau+1}(\mathbf{V}_f)$ and $\mathbf{C}_{\tau+1}(\mathbf{V}_c)$ are all very close to $\mathbf{0}$. Hence, if an output cell error occurs only at the end of the sequence, *i.e.* $\boldsymbol{\lambda}_{t,t}$ is significantly different from zero, such a memory cell, via backpropagation through time, supplies the same cell error rate $\boldsymbol{\lambda}_{t,\tau}$ in every timestep $\tau = 1, \dots, t$ that would not vanish and stabilize. Thus, all weights \mathbf{W} 's, \mathbf{V} 's, and biases \mathbf{b} 's of a memory cell are updated by the same cell error rate $\boldsymbol{\lambda}_{t,t}$, and thus, LSTM avoids the vanishing gradient problem through time.

Moreover, with the input gate \mathbf{i} being close to zero, the memory cell just sums up all the fed-in inputs it receives during scanning the input temporal sequence, *i.e.* $\mathbf{i}_\tau = \mathbf{0}$, for all $\tau = 1, \dots, t$,

$$\mathbf{c}_\tau = \mathbf{i}_\tau \odot \tilde{\mathbf{c}}_\tau + \mathbf{f}_\tau \odot \mathbf{c}_{\tau-1} \quad \mapsto \quad \mathbf{c}_\tau = \tilde{\mathbf{c}}_\tau + \mathbf{f}_\tau \odot \mathbf{c}_{\tau-1}.$$

Thus, such a memory cell is the same as a unit that reads in all sequence elements at one time.

Bibliography

- [1] Burkov, A. (2019). *The hundred-page machine learning book*, volume 1. Andriy Burkov Quebec City, QC, Canada.
- [2] Cho, K., Van Merriënboer, B., Bahdanau, D., and Bengio, Y. (2014). On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*.
- [3] Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.
- [4] Olah, C. (2015). Understanding lstm networks.