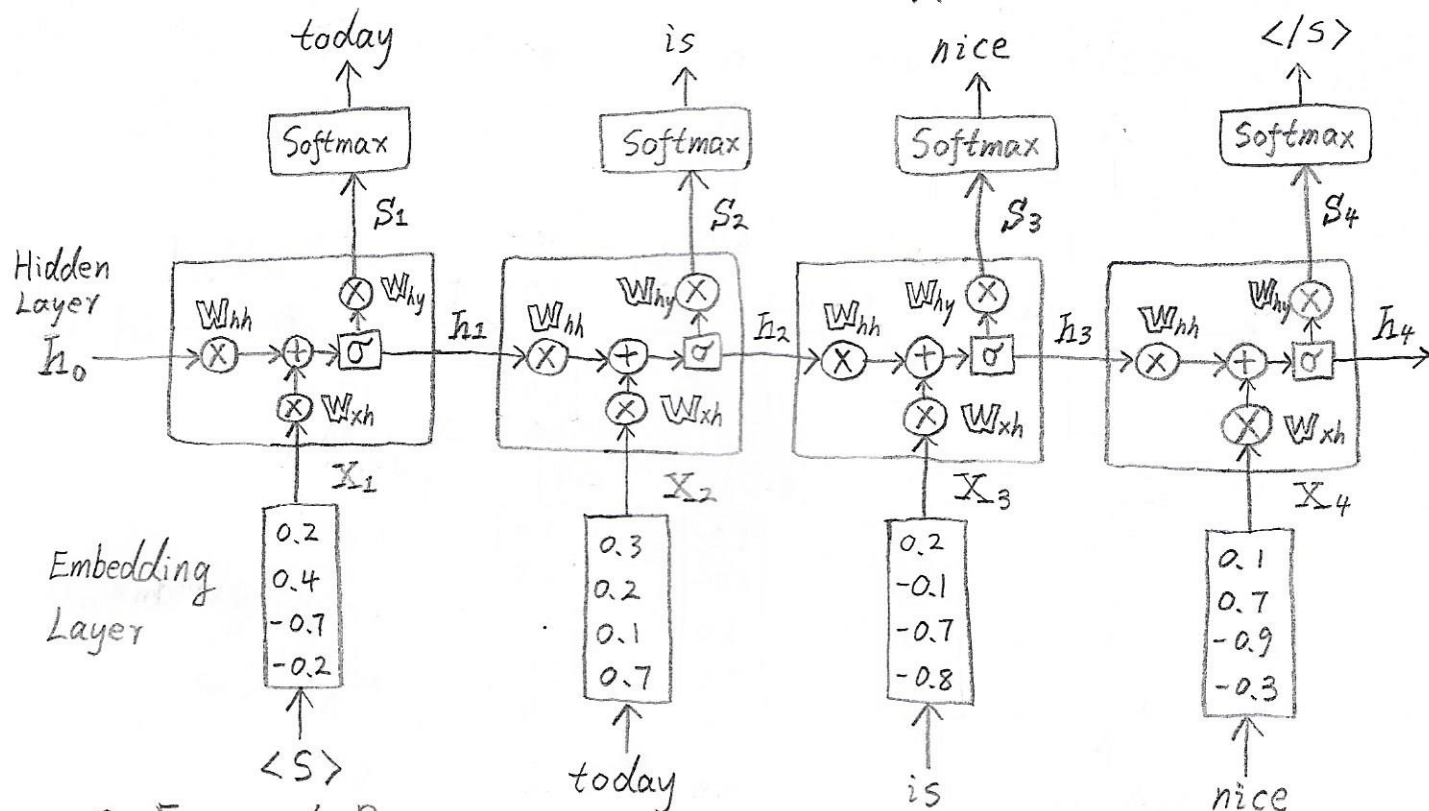# IX. RECURRENT NEURAL NETWORK

## 1. Vanilla Recurrent Neural Network



## 2. Forward Pass

$$W_e = \left\{ \text{"hello"}: X_1, \text{"what"}: X_2, \dots, \text{"zaltan"}: X_{v\text{-size}} \right\}$$

$$X_t = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{256} \end{bmatrix}, \quad \text{hidden\_size} = 256$$

$$\underset{256 \times 1}{h_t} = \sigma( \underset{256 \times 256}{W_{xh}} \cdot \underset{256 \times 1}{X_t} + \underset{256 \times 256}{W_{hh}} \cdot \underset{256 \times 1}{h_{t-1}})$$

$$= \frac{1}{1 + e^{-(W_{xh} \cdot X_t + W_{hh} \cdot h_{t-1})}}$$

$$\underset{|Y| \times 1}{S_t} = \underset{|Y| \times 256}{W_{hy}} \cdot \underset{256 \times 1}{h_t} = \begin{bmatrix} S_t(y_1) \\ S_t(y_2) \\ \vdots \\ S_t(y_m) \end{bmatrix}, \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

193

$$P_t = \text{softmax}(S_t)$$

$$= \text{softmax}\left( \begin{bmatrix} S_t(y_1) \\ S_t(y_2) \\ \vdots \\ S_t(y_m) \end{bmatrix} \right)$$

$$= \begin{bmatrix} \dfrac{e^{S_t(y_1)}}{\sum_{y' \in Y} e^{S_t(y')}} \\ \vdots \\ \dfrac{e^{S_t(y_m)}}{\sum_{y' \in Y} e^{S_t(y')}} \end{bmatrix} = \begin{bmatrix} P_t(y_1) \\ \vdots \\ P_t(y_m) \end{bmatrix} = \begin{bmatrix} P_t(y_1) \\ \vdots \\ P_t(y_t) \\ \vdots \\ P_t(y_m) \end{bmatrix}$$

True output at index $i$ → Estimated output at index $i$.

$$loss_t = -\sum_{i=1}^{m} y_i' \cdot \log(P_t(y_i)) = -y_t' \cdot \log(P_t(y_t))$$

( If $y_i' = y_t'$
$y_i' = 1$;
if $y_i' \neq y_t'$,
$y_i' = 0$.
Cross Entropy
Loss )

$$= -\log\left( \frac{e^{S_t(y_t)}}{\sum_{y' \in Y} e^{S_t(y')}} \right)$$

$$= \log \sum_{y' \in Y} e^{S_t(y')} - \log e^{S_t(y_t)}$$

$$= \log \sum_{y' \in Y} e^{S_t(y')} - S_t(y_t)$$

3. Backward Pass

$$dS_t(y) = \frac{\partial loss_t}{\partial S_t(y)} = \frac{\partial}{\partial S_t(y)}\left( \log \sum_{y' \in Y} e^{S_t(y')} - S_t(y_t) \right)$$

$$= \frac{\partial}{\partial S_t(y)}\left( \log \sum_{y' \in Y} e^{S_t(y')} \right) - \frac{\partial}{\partial S_t(y)}(S_t(y_t))$$

$$= \frac{1}{\sum_{y' \in Y} e^{S_t(y')}} \cdot \frac{\partial}{\partial S_t(y)}\left( \sum_{y' \in Y} e^{S_t(y')} \right) - \frac{\partial}{\partial S_t(y)}(S_t(y_t))$$

⟨ chain rule ⟩

When $y = y_t$

$$d S_t(y_t) = \frac{e^{S_t(y_t)}}{\sum\limits_{y' \in Y} e^{S_t(y')}} - 1$$

$$= P_t(y_t) - 1$$

When $y \neq y_t$

$$d S_t(y) = \frac{e^{S_t(y)}}{\sum\limits_{y' \in Y} e^{S_t(y')}} = P_t(y)$$

$$d S_t = \frac{\partial \text{loss}_t}{\partial S_t} = \begin{bmatrix} d S_t(y_1) \\ \vdots \\ d S_t(y_t) \\ \vdots \\ d S_t(y_m) \end{bmatrix} = \begin{bmatrix} P_t(y_1) \\ \vdots \\ P_t(y_t) - 1 \\ \vdots \\ P_t(y_m) \end{bmatrix}$$

$$= \begin{bmatrix} P_t(y_1) \\ \vdots \\ P_t(y_t) \\ \vdots \\ P_t(y_m) \end{bmatrix} - \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} = P_t - 1_{y = y_t} \quad \circledcirc$$

# 4. Mathematical Helpers

## 1. Definition 1

If $U$ and $V$ are vectors, then

$$\text{diag}(U) \cdot V$$

$$= \text{diag}\left(\begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}\right) \cdot \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} u_1 & 0 & \cdots & 0 \\ 0 & u_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & u_n \end{bmatrix} \cdot \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} u_1 v_1 \\ u_2 v_2 \\ \vdots \\ u_n v_n \end{bmatrix} = U \circ V$$

## 2. Lemma 1

$V$ and $h$ are vectors. $W$ is a matrix. $dV = \dfrac{\partial l}{\partial V}$ is known.

$V = f(W \cdot h)$, where $f$ is an element-wise function. Then

$$dh = \frac{\partial l}{\partial h} = W^T \cdot (f'(W \cdot h) \circ dV), \quad \text{①} \quad \text{and}$$

$$dW = \frac{\partial l}{\partial W} = (f'(W \cdot h) \circ dV) \cdot h^T \quad \text{②}$$

Proof:

$$l = g(f(W \cdot h)), \quad \text{where} \quad V = f(W \cdot h), \quad Z = W \cdot h$$

$$\frac{\partial l}{\partial h} = \frac{\partial f}{\partial h} \cdot \frac{\partial l}{\partial f} \quad (\text{Chain rule, D.L., P104, Denominator-Layout})$$

$$= \left(\frac{\partial Z}{\partial h} \cdot \frac{\partial f}{\partial Z}\right) \cdot \frac{\partial l}{\partial f} \quad (\text{Chain rule, again})$$

$$= \frac{\partial Wh}{\partial h} \cdot \frac{\partial f(Z)}{\partial Z} \cdot \frac{\partial l}{\partial V}$$

$$= \frac{\partial Wh}{\partial h} \cdot \text{diag}(f'(Z)) \cdot dV$$

$$= \frac{\partial \mathbb{W} \cdot h}{\partial h} \cdot (f'(\mathbb{W}h) \circ dv) \quad (\text{Definition 1})$$

$$\mathbb{W} = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1n} \\ w_{21} & w_{22} & \cdots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1} & w_{m2} & \cdots & w_{mn} \end{bmatrix}, \quad h = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_n \end{bmatrix}, \quad \text{then}$$

$$\mathbb{W}h = \begin{bmatrix} w_{11}h_1 + w_{12}h_2 + \cdots + w_{1n}h_n \\ w_{21}h_1 + w_{22}h_2 + \cdots + w_{2n}h_n \\ \vdots \\ w_{m1}h_1 + w_{m2}h_2 + \cdots + w_{mn}h_n \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix}$$

$$\frac{\partial \mathbb{W}h}{\partial h} = \begin{bmatrix} \frac{\partial z_1}{\partial h_1} & \frac{\partial z_2}{\partial h_1} & \cdots & \frac{\partial z_m}{\partial h_1} \\ \frac{\partial z_1}{\partial h_2} & \frac{\partial z_2}{\partial h_2} & \cdots & \frac{\partial z_m}{\partial h_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_1}{\partial h_n} & \frac{\partial z_2}{\partial h_n} & \cdots & \frac{\partial z_m}{\partial h_n} \end{bmatrix} = \begin{bmatrix} w_{11} & w_{21} & \cdots & w_{m1} \\ w_{12} & w_{22} & \cdots & w_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ w_{1n} & w_{2n} & \cdots & w_{mn} \end{bmatrix} = \mathbb{W}^T$$

$$(\text{Denominator - Layout})$$

$$dh = \frac{\partial \ell}{\partial h} = \mathbb{W}^T \cdot (f'(\mathbb{W}h) \circ dv)$$

$$\text{Let} \quad w_1 = \begin{bmatrix} w_{11} \\ w_{12} \\ \vdots \\ w_{1n} \end{bmatrix}, \quad w_2 = \begin{bmatrix} w_{21} \\ w_{22} \\ \vdots \\ w_{2n} \end{bmatrix}, \quad \cdots, \quad w_m = \begin{bmatrix} w_{m1} \\ w_{m2} \\ \vdots \\ w_{mn} \end{bmatrix}, \quad \text{then}$$

$$\ell = g(f(\mathbb{W}h))$$

$$= g\left( f\left( \begin{bmatrix} w_1^T \\ w_2^T \\ \vdots \\ w_m^T \end{bmatrix} \cdot h \right) \right)$$

$$= g\left(f\left(\begin{bmatrix} W_1^T \cdot h \\ W_2^T \cdot h \\ \vdots \\ W_m^T \cdot h \end{bmatrix}\right)\right) = g\left(f\left(\begin{bmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_m \end{bmatrix}\right)\right)$$

$$= g\left(\begin{bmatrix} f(Z_1) \\ f(Z_2) \\ \vdots \\ f(Z_m) \end{bmatrix}\right) = g\left(\begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_m \end{bmatrix}\right)$$

$$\frac{\partial \ell}{\partial W_i} = \frac{\partial v_i}{\partial W_i} \cdot \frac{\partial \ell}{\partial v_i} \quad (\text{Chain Rule})$$

$$= \left(\frac{\partial z_i}{\partial W_i} \cdot \frac{\partial v_i}{\partial z_i}\right) \cdot \frac{\partial \ell}{\partial v_i} \quad (\text{Chain Rule})$$

$$= \frac{\partial W_i^T h}{\partial W_i} \cdot f'(Z_i) \cdot dv_i$$

$$= \begin{bmatrix} \dfrac{\partial(w_{i1} h_1 + w_{i2} h_2 + \cdots + w_{in} h_n)}{\partial w_{i1}} \\ \vdots \\ \dfrac{\partial(w_{i1} h_1 + w_{i2} h_2 + \cdots + w_{in} h_n)}{\partial w_{in}} \end{bmatrix} \cdot f'(Z_i) \cdot dv_i$$

$$= \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_n \end{bmatrix} \cdot f'(Z_i) \cdot dv_i$$

$$= h \cdot f'(z_i) \cdot dv_i$$

$$dw_i^T = \frac{\partial \ell}{\partial w_i^T} = f'(z_i) \cdot dv_i \cdot h^T$$

$$dW = \frac{\partial \ell}{\partial W} = \begin{bmatrix} \frac{\partial \ell}{\partial w_1^T} \\ \frac{\partial \ell}{\partial w_2^T} \\ \vdots \\ \frac{\partial \ell}{\partial w_m^T} \end{bmatrix} = \begin{bmatrix} f'(z_1) \cdot dv_1 \\ f'(z_2) \cdot dv_2 \\ \vdots \\ f'(z_m) \cdot dv_m \end{bmatrix} \cdot h^T$$

$$= (f'(W \cdot h) \circ dv) \cdot h^T \qquad \square$$

③. Corollary 1

If $v = f(Wh) = Wh$, then

$$dh = W^T \cdot dv \qquad ③$$

$$dW = dv \cdot h^T \qquad ④$$

Proof:

$$f(Wh) = f\left(\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix}\right) = \begin{bmatrix} f(z_1) \\ f(z_2) \\ \vdots \\ f(z_m) \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix}$$

$$f'(Wh) = f'\left(\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix}\right) = \begin{bmatrix} f'(z_1) \\ f'(z_2) \\ \vdots \\ f'(z_m) \end{bmatrix} = \begin{bmatrix} \frac{dz_1}{dz_1} \\ \frac{dz_2}{dz_2} \\ \vdots \\ \frac{dz_m}{dz_m} \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$\therefore dh = W^T \cdot dv, \quad dW = dv \cdot h^T \quad (\text{From } ① \text{ and } ②) \qquad \square$$

4) Lemma 2

$U$, $V$ and $S$ are vectors. $S = U \circ f(V)$.

$$dU = \frac{\partial l}{\partial U}, \quad dV = \frac{\partial l}{\partial V}, \quad dS = \frac{\partial l}{\partial S}, \quad \text{then}$$

$$dU = f(V) \circ dS \qquad \text{⑤}$$

$$dV = f'(V) \circ U \circ dS \qquad \text{⑥}$$

Proof:

$$l = g(S) = g(U \circ f(V))$$

$$dU = \frac{\partial S}{\partial U} \cdot \frac{\partial l}{\partial S} = f(V) \circ dS$$

$$dV = \frac{\partial S}{\partial V} \cdot \frac{\partial l}{\partial S} = \frac{\partial f}{\partial V} \cdot \frac{\partial S}{\partial f} \cdot \frac{\partial l}{\partial S} = f'(V) \circ U \circ dS \qquad \square$$

5) Corollary 2

If $f(V) = V$, then

$$dU = V \circ dS \qquad \text{⑦}$$

$$dV = U \circ dS \qquad \text{⑧}$$

5. Backward Pass – Continue

Since $S_t = W_{hy} \cdot h_t$,

$$dh_t = W_{hy}^T \cdot dS_t \quad \langle by\ ③\rangle. \quad ⑨$$

$$d W_{hy} = dS_t \cdot h_t^T \quad \langle by\ ④\rangle \quad ⑩$$

$$W_{xh} \cdot X_t + W_{hh} \cdot h_{t-1}$$

$$= \begin{bmatrix} W_{xh11} & W_{xh12} & \cdots & W_{xh1n} \\ W_{xh21} & W_{xh22} & \cdots & W_{xh2n} \\ \vdots & \vdots & \ddots & \vdots \\ W_{xhm1} & W_{xhm2} & \cdots & W_{xhmn} \end{bmatrix} \cdot \begin{bmatrix} X_{t1} \\ X_{t2} \\ \vdots \\ X_{tn} \end{bmatrix} + \begin{bmatrix} W_{hh11} & W_{hh12} & \cdots & W_{hh1n} \\ W_{hh21} & W_{hh22} & \cdots & W_{hh2n} \\ \vdots & \vdots & \ddots & \vdots \\ W_{hhm1} & W_{hhm2} & \cdots & W_{hhmn} \end{bmatrix} \cdot \begin{bmatrix} h_{t-11} \\ h_{t-12} \\ \vdots \\ h_{t-1n} \end{bmatrix}$$

$$= \begin{bmatrix} W_{xh11} & W_{xh12} & \cdots & W_{xh1n} & W_{hh11} & W_{hh12} & \cdots & W_{hh1n} \\ W_{xh21} & W_{xh22} & \cdots & W_{xh2n} & W_{hh21} & W_{hh22} & \cdots & W_{hh2n} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ W_{xhm1} & W_{xhm2} & \cdots & W_{xhmn} & W_{hhm1} & W_{hhm2} & \cdots & W_{hhmn} \end{bmatrix} \cdot \begin{bmatrix} X_{t1} \\ X_{t2} \\ \vdots \\ X_{tn} \\ h_{t-11} \\ h_{t-12} \\ \vdots \\ h_{t-1n} \end{bmatrix}$$

$$= \begin{bmatrix} W_{xh} & W_{hh} \end{bmatrix} \cdot \begin{bmatrix} X_t \\ h_{t-1} \end{bmatrix}$$

$$= \begin{bmatrix} W_{xh} & W_{hh} \end{bmatrix} \cdot \begin{bmatrix} X_t ; h_{t-1} \end{bmatrix}$$

$$= T_{rnn} \cdot Z_t$$

$$\therefore \quad h_t = \sigma(T_{rnn} \cdot Z_t)$$

$$dz_t = T_{rnn}^T \cdot (\sigma'(T_{rnn} Z_t) \circ dh_t) \quad \langle by \ ① \rangle \quad ⑪$$

$$\begin{bmatrix} dX_t \\ dh_{t-1} \end{bmatrix} = \begin{bmatrix} W_{xh}^T \\ W_{hh}^T \end{bmatrix} \cdot (\sigma'(T_{rnn} Z_t) \circ dh_t)$$

$$dX_t = W_{xh}^T \cdot (\sigma'(T_{rnn} Z_t) \circ dh_t) \quad ⑫$$

$$dh_{t-1} = W_{hh}^T \cdot (\sigma'(T_{rnn} Z_t) \circ dh_t) \quad ⑬$$

$$dT_{rnn} = (\sigma'(T_{rnn} Z_t) \circ dh_t) \cdot Z_t^T \quad \langle by \ ② \rangle \quad ⑭$$

$$\begin{bmatrix} dW_{xh} & dW_{hh} \end{bmatrix} = (\sigma'(T_{rnn} Z_t) \circ dh_t) \cdot \begin{bmatrix} X_t^T & h_{t-1}^T \end{bmatrix}$$

$$dW_{xh} = (\sigma'(T_{rnn} Z_t) \circ dh_t) \cdot X_t^T \quad ⑮$$

$$dW_{hh} = (\sigma'(T_{rnn} Z_t) \circ dh_t) \cdot h_{t-1}^T \quad ⑯$$

Backpropagation_Through_Time_Vanilla_RNN ( ) :

for t from T to 1 :

$$\frac{\partial l_t}{\partial S_t} = P_t - 1 \; y = y_t \quad ⓪$$

$$\frac{\partial l_t}{\partial W_{hy}} = \frac{\partial l_{t+1}}{\partial W_{hy}} + \frac{\partial l_t}{\partial S_t} \cdot h^T \quad ⑩$$

$$\frac{\partial l_t}{\partial h_t} = \frac{\partial l_{t+1}}{\partial h_t} + W_{hy}^T \cdot \frac{\partial l_t}{\partial S_t} \quad ⑨$$

$$\frac{\partial l_t}{\partial W_{xh}} = \frac{\partial l_{t+1}}{\partial W_{xh}} + \left( \sigma'(T_{rnn} Z_t) \circ \frac{\partial l_t}{\partial h_t} \right) \cdot X_t^T \quad ⑮$$

$$\frac{\partial l_t}{\partial W_{hh}} = \frac{\partial l_{t+1}}{\partial W_{hh}} + \left( \sigma'(T_{rnn} Z_t) \circ \frac{\partial l_t}{\partial h_t} \right) \cdot h_{t-1}^T \quad ⑯$$

$$\frac{\partial l_t}{\partial X_t} = W_{xh}^T \cdot \left( \sigma'(T_{rnn} Z_t) \circ \frac{\partial l_t}{\partial h_t} \right) \quad ⑫$$

$$\frac{\partial l_t}{\partial h_{t-1}} = W_{hh}^T \cdot \left( \sigma'(T_{rnn} Z_t) \circ \frac{\partial l_t}{\partial h_t} \right) \quad ⑬$$

Comments :

1) $\frac{\partial l_t}{\partial W_{hy}}$, $\frac{\partial l_t}{\partial W_{xh}}$, $\frac{\partial l_t}{\partial W_{hh}}$, $\frac{\partial l_t}{\partial h_t}$ are aggregated through time since

$$\frac{\partial l_{total}}{\partial W} = \frac{\partial l_{t=T} + l_{t=T-1} + \cdots + l_{t=1}}{\partial W} = \frac{\partial l_{t=T}}{\partial W} + \frac{\partial l_{t=T-1}}{\partial W} + \cdots + \frac{\partial l_{t=1}}{\partial W}$$

2) $\sigma'(T_{rnn} Z_t) = \sigma(T_{rnn} Z_t) - \sigma^2(T_{rnn} Z_t) = \sigma\left( [W_{xh} \; W_{hh}] \cdot \begin{bmatrix} X_t \\ h_{t-1} \end{bmatrix} \right) - \sigma^2\left( [W_{xh} \; W_{hh}] \begin{bmatrix} X_t \\ h_{t} \end{bmatrix} \right)$

( D.L., P 48 )

203