

---

# A Novel Optimization Method for Recurrent Neural Network by Non-Euclidean Geometry - Appendix

---

<b>Chen Liu</b> Master Student, IC, EPFL chen.liu@epfl.ch	Ya-Ping Hsieh Supervisor, PH.D, EL, EPFL ya-ping.hsieh@epfl.ch
Volkan Cevher Supervisor, Associate Professor, EL, EPFL volkan.cevher@epfl.ch	

## A General Upper Bound of RNN Loss Function

Let  $\Phi = \{\mathbf{U}, \mathbf{W}, \mathbf{V}, \mathbf{s}_0\}$ . We already have the expression of t-th token's contribution to loss function.

$$E_t(\Phi) = \log \sum_{j=1}^K e^{\mathbf{V}_j \mathbf{s}_t} - \mathbf{y}^T \mathbf{V} \mathbf{s}_t \quad (1)$$

This function contains two parts: log-sum-exp part and linear part. Let  $f_t(\Phi) = \log \sum_{j=1}^K e^{\mathbf{V}_j \mathbf{s}_t}$ ,  $g_t(\Phi) = -\mathbf{y}^T \mathbf{V} \mathbf{s}_t$  and  $\mathbf{r} = \mathbf{V} \mathbf{s}_t(\Phi + \Delta \Phi) - \mathbf{V} \mathbf{s}_t(\Phi) - \langle \nabla_{\Phi} \mathbf{V} \mathbf{s}_t(\Phi), \Delta \Phi \rangle$ . We will obtain the upper bound of formula 1 using the following two lemmas.

**Lemma 1**  $f_t(\Phi + \Delta \Phi) \leq f_t(\Phi) + \langle \nabla_{\Phi} f_t(\Phi), \Delta \Phi \rangle + \frac{1}{2} \|\mathbf{V} \mathbf{s}_t(\Phi + \Delta \Phi) - \mathbf{V} \mathbf{s}_t(\Phi)\|_{\infty}^2 + \|\mathbf{r}\|_{\infty}$

According to [1], for log-sum-exp function  $h(\mathbf{x}) = \log \sum_{j=1}^J e^{\mathbf{x}_j}$ , Lipschitz constant  $L_{\infty} = \frac{1}{2}$  and we can derive an upper bound of function  $h(\mathbf{x})$ .

$$h(\mathbf{x} + \Delta \mathbf{x}) \leq h(\mathbf{x}) + \langle \nabla_{\mathbf{x}} h(\mathbf{x}), \Delta \mathbf{x} \rangle + \frac{1}{2} \|\Delta \mathbf{x}\|_{\infty}^2 \quad (2)$$

It is not difficult to find  $f_t(\Phi) = h(\mathbf{V} \mathbf{s}_t(\Phi))$ . Plug this into equation 2, we have the following formula.

$$f_t(\Phi + \Delta \Phi) \leq f_t(\Phi) + \langle \nabla_{\mathbf{V} \mathbf{s}_t} f_t(\Phi), \mathbf{V} \mathbf{s}_t(\Phi + \Delta \Phi) - \mathbf{V} \mathbf{s}_t(\Phi) \rangle + \frac{1}{2} \|\mathbf{V} \mathbf{s}_t(\Phi + \Delta \Phi) - \mathbf{V} \mathbf{s}_t(\Phi)\|_{\infty}^2 \quad (3)$$

From the definition of  $\mathbf{r}$ , we have  $\mathbf{V} \mathbf{s}_t(\Phi + \Delta \Phi) - \mathbf{V} \mathbf{s}_t(\Phi) = \mathbf{r} + \langle \nabla_{\Phi} \mathbf{V} \mathbf{s}_t(\Phi), \Delta \Phi \rangle$ . Subsequently, we have the following equation.

$$\begin{aligned} \langle \nabla_{\mathbf{V} \mathbf{s}_t} f_t(\Phi), \mathbf{V} \mathbf{s}_t(\Phi + \Delta \Phi) - \mathbf{V} \mathbf{s}_t(\Phi) \rangle &= \langle \nabla_{\mathbf{V} \mathbf{s}_t} f_t(\Phi), \mathbf{r} + \langle \nabla_{\Phi} \mathbf{V} \mathbf{s}_t(\Phi), \Delta \Phi \rangle \rangle \\ &= \mathbf{r}^T \nabla_{\mathbf{V} \mathbf{s}_t} f_t(\Phi) + \langle \nabla_{\mathbf{V} \mathbf{s}_t} f_t(\Phi), \nabla_{\Phi} \mathbf{V} \mathbf{s}_t(\Phi), \Delta \Phi \rangle = \mathbf{r}^T \nabla_{\mathbf{V} \mathbf{s}_t} f_t(\Phi) + \langle \nabla_{\Phi} f_t(\Phi), \Delta \Phi \rangle \end{aligned} \quad (4)$$

From  $\nabla_{\mathbf{V} \mathbf{s}_t} f_t(\Phi) = \text{softmax}(\mathbf{V} \mathbf{s}_t)$ , we have  $\mathbf{r}^T \nabla_{\mathbf{V} \mathbf{s}_t} f_t(\Phi) \leq \|\mathbf{r}\|_{\infty} \|\text{softmax}(\mathbf{V} \mathbf{s}_t)\|_{\infty} \leq \|\mathbf{r}\|_{\infty}$ . Combine this with equation 3 and 4, we will get **Lemma 1**. #

**Lemma 2**  $g_t(\Phi + \Delta \Phi) \leq \langle \nabla_{\Phi} g_t(\Phi), \Delta \Phi \rangle + \|\mathbf{r}\|_{\infty}$

$g_t(\Phi)$  is linear, so it is equal to its first-order Taylor expansion.

$$g_t(\Phi + \Delta \Phi) = g_t(\Phi) + \langle \nabla_{\mathbf{V} \mathbf{s}_t} g_t(\Phi), \mathbf{V} \mathbf{s}_t(\Phi + \Delta \Phi) - \mathbf{V} \mathbf{s}_t(\Phi) \rangle \quad (5)$$

Similar to formula 4, we have the following upper bound expression.

$$\langle \nabla_{\mathbf{V}_{\mathbf{s}_t} g_t(\Phi)}, \mathbf{V}_{\mathbf{s}_t}(\Phi + \Delta\Phi) - \mathbf{V}_{\mathbf{s}_t}(\Phi) \rangle \leq \langle \nabla_{\Phi} g_t(\Phi), \Delta\Phi \rangle + \|\mathbf{r}\|_{\infty} \quad (6)$$

Combine formula 5 with 6, we will get **Lemma 2** immediately. #

The general upper bound of RNN loss function below is obtained by adding **Lemma 1** and **Lemma 2**.

$$\begin{aligned} E_t(\Phi + \Delta\Phi) &\leq E_t(\Phi) + \langle \nabla_{\Phi} E_t(\Phi), \Delta\Phi \rangle + \frac{1}{2} \|\mathbf{V}_{\mathbf{s}_t}(\Phi + \Delta\Phi) - \mathbf{V}_{\mathbf{s}_t}(\Phi)\|_{\infty}^2 + \\ &\quad 2 \|\mathbf{V}_{\mathbf{s}_t}(\Phi + \Delta\Phi) - \mathbf{V}_{\mathbf{s}_t}(\Phi) - \langle \mathbf{V}_{\nabla \mathbf{s}_t}(\Phi), \Delta\Phi \rangle\|_{\infty} \end{aligned} \quad (7)$$

## B Lipschitz Constants of Some Matrices

In this section, we will derive the upper bound of Lipschitz constants of RNN's input matrix and recurrent matrix. We will repeatedly use the following inequality in the following subsections. The proof of them are available in the section 9.3 in [2] or in the appendix of [3].

$$\text{if } q^{-1} + p^{-1} = r^{-1} \text{ and } p, q \in [1, +\infty], \text{ then } \|\mathbf{M}\|_{S^r} \leq \|\mathbf{M}\|_{S^q} \|\mathbf{M}\|_{S^p} \quad (8)$$

$$\text{if } \mathbf{v} \text{ is a vector, then } \|\mathbf{r}\|_{S^p} = \|\mathbf{r}\|_2 \quad \forall p \in [1, +\infty] \quad (9)$$

$$\text{tr}(\mathbf{M}_1 \mathbf{M}_2) \leq \|\mathbf{M}_1\|_{S^1} \|\mathbf{M}_2\|_{S^{\infty}} \quad (10)$$

$$\|\mathbf{x} \odot \mathbf{y}\|_2 \leq \|\mathbf{x}\|_{\infty} \|\mathbf{y}\|_2 \quad (11)$$

Here,  $\|\mathbf{M}\|_{S^p}$  means Schatten- $p$  norm of a matrix,  $\text{tr}(\mathbf{M})$  is the trace of a matrix and  $\odot$  means element-wise multiplication. The notation in this section is the same as the main body.

### B.1 Lipschitz Constant of the Input Matrix

According to inequality 7, we have two parts to be upper bounded:  $\|\mathbf{V}_{\mathbf{s}_t}(\mathbf{U}^k + \Delta\mathbf{U}) - \mathbf{V}_{\mathbf{s}_t}(\mathbf{U}^k)\|_{\infty}$  and  $\|\mathbf{V}_{\mathbf{s}_t}(\mathbf{U}^k + \Delta\mathbf{U}) - \mathbf{V}_{\mathbf{s}_t}(\mathbf{U}^k) - \langle \mathbf{V}_{\nabla \mathbf{s}_t}(\mathbf{U}^k), \Delta\mathbf{U} \rangle\|_{\infty}$ .

The first part is easy. We can approximate  $\mathbf{s}_t(\mathbf{U} + \Delta\mathbf{U})$  using its first order Taylor expansion.

$$[\mathbf{V}_{\mathbf{s}_t}(\mathbf{U}^k + \Delta\mathbf{U}) - \mathbf{V}_{\mathbf{s}_t}(\mathbf{U}^k)]_p \simeq \text{tr} \left[ \frac{\partial [\mathbf{V}_{\mathbf{s}_t}]_p}{\partial \mathbf{U}} \Delta\mathbf{U} \right] \leq \|\mathbf{Q}_t \tilde{\mathbf{V}}_p\|_{S^1} \|\Delta\mathbf{U}\|_{S^{\infty}} \quad (12)$$

The last less than sign satisfies using formula 10 and  $\frac{\partial \mathbf{s}_t(\mathbf{U})}{\partial \mathbf{U}} = \mathbf{Q}_t$ . So, the upper bound of first part is:

$$\|\mathbf{V}_{\mathbf{s}_t}(\mathbf{U}^k + \Delta\mathbf{U}) - \mathbf{V}_{\mathbf{s}_t}(\mathbf{U}^k)\|_{\infty} \leq \max_p \|\mathbf{Q}_t \tilde{\mathbf{V}}_p\|_{S^1} \|\Delta\mathbf{U}\|_{S^{\infty}} \quad (13)$$

For the second part, we use Taylor expansion to analyze its  $p$ -th element.

$$\begin{aligned} &[\mathbf{V}_{\mathbf{s}_t}(\mathbf{U}^k + \Delta\mathbf{U}) - \mathbf{V}_{\mathbf{s}_t}(\mathbf{U}^k) - \langle \mathbf{V}_{\nabla \mathbf{U} \mathbf{s}_t}(\mathbf{U}^k), \Delta\mathbf{U} \rangle]_p \\ &= \int_0^1 \text{tr}([\mathbf{V}_{p,:} \nabla_{\mathbf{U} \mathbf{s}_t}(\mathbf{U}^k + t\Delta\mathbf{U}) - \mathbf{V}_{p,:} \nabla_{\mathbf{U} \mathbf{s}_t}(\mathbf{U}^k)] \Delta\mathbf{U}) dt \\ &\leq \|\Delta\mathbf{U}\|_{S^{\infty}} \int_0^1 \|\mathbf{V}_{p,:} \nabla_{\mathbf{U} \mathbf{s}_t}(\mathbf{U}^k + t\Delta\mathbf{U}) - \mathbf{V}_{p,:} \nabla_{\mathbf{U} \mathbf{s}_t}(\mathbf{U}^k)\|_{S^1} dt \end{aligned} \quad (14)$$

We focus on the integration part and like first part we use first-order Taylor expansion approximation.

$$\begin{aligned}
& \int_0^1 \mathbf{V}_{p,:} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k + t\Delta\mathbf{U}) - \mathbf{V}_{p,:} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k) \\
& \simeq \int_0^1 \frac{d}{dt} \mathbf{V}_{p,:} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k + t\Delta\mathbf{U})|_{t=0} t dt \\
& = \frac{1}{2} \frac{d}{dt} \mathbf{V}_{p,:} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k + t\Delta\mathbf{U})|_{t=0}
\end{aligned} \tag{15}$$

From RNN's hidden states iterative equation  $\mathbf{s}_t = \sigma(\mathbf{W}\mathbf{s}_{t-1} + \mathbf{U}\mathbf{x}_t)$ , we can calculate  $\mathbf{U}$ 's first order gradient:  $\nabla_{\mathbf{U}} \mathbf{V}_p \mathbf{s}_t = \Lambda'_t \mathbf{V}_p^T \mathbf{x}_t + \mathbf{Q}_{t-1} \tilde{\mathbf{W}}^T \tilde{\Lambda}'_t \tilde{\mathbf{V}}_p^T$  where  $\Lambda'_t$  is a diagonal matrix whose diagonal vector is  $\sigma'(\mathbf{W}\mathbf{s}_{t-1} + \mathbf{U}\mathbf{x}_t)$ . Subsequently, we can obtain the following equation where we ignore the second order part  $\frac{d\mathbf{Q}_{t-1}}{dt} = tr(\frac{\partial^2 \mathbf{s}_{t-1}}{\partial \mathbf{U}^2} \Delta\mathbf{U})$ .

$$\frac{d}{dt} \nabla_{\mathbf{U}} \mathbf{V}_p \mathbf{s}_t = \frac{d}{dt} \Lambda'_t \mathbf{V}_p^T \mathbf{x}_t + \mathbf{Q}_{t-1} \tilde{\mathbf{W}}^T \frac{d}{dt} \tilde{\Lambda}'_t \tilde{\mathbf{V}}_p^T \tag{16}$$

Then we can obtain the upper bound containing  $\|\mathbf{U}\|_{S^\infty}$  of the equation above using equation 8, 9, and 11.

$$\begin{aligned}
\|\frac{d}{dt} \Lambda'_t \mathbf{V}_p^T \mathbf{x}_t\|_{S^1} & \leq \|\frac{d}{dt} \lambda'_t \odot \mathbf{V}_p^T\|_2 \|\mathbf{x}_t\|_2 \leq \|\frac{d}{dt} \lambda'_t\|_\infty \|\mathbf{V}_p\|_2 \|\mathbf{x}_t\|_2 \\
& \leq \max_q (\lambda'^{-1}_{t_q} \lambda''_{t_q} \|\mathbf{Q}_{t-1, :, q}\|_{S^1}) \|\Delta\mathbf{U}\|_{S^\infty} \|\mathbf{V}_p\|_2 \|\mathbf{x}_t\|_2
\end{aligned} \tag{17}$$

$$\begin{aligned}
& \|\mathbf{Q}_{t-1} \tilde{\mathbf{W}}^T \frac{d}{dt} \tilde{\Lambda}'_t \tilde{\mathbf{V}}_p^T\|_{S^1} \leq \|\mathbf{Q}_{t-1, :, *}\|_{S^1} \|\mathbf{W}^T \frac{d}{dt} \lambda'_t \odot \mathbf{V}_p^T\|_2 \\
& \leq \|\mathbf{Q}_{t-1, :, *}\|_{S^1} \|\mathbf{W}\|_{S^\infty} \|\frac{d}{dt} \lambda'_t \odot \mathbf{V}_p^T\|_2 \leq \|\mathbf{Q}_{t-1, :, *}\|_{S^1} \|\mathbf{W}\|_{S^\infty} \|\frac{d}{dt} \lambda'_t\|_\infty \|\mathbf{V}_p\|_2 \\
& \leq \|\mathbf{Q}_{t-1, :, *}\|_{S^1} \|\mathbf{W}\|_{S^\infty} \max_q (\lambda'^{-1}_{t_q} \lambda''_{t_q} \|\mathbf{Q}_{t-1, :, q}\|_{S^1}) \|\Delta\mathbf{U}\|_{S^\infty} \|\mathbf{V}_p\|_2
\end{aligned} \tag{18}$$

The last less than sign of both formulas and the first less than sign of the second one use the following two lemmas respectively.

**Lemma 3**  $\|\frac{d}{dt} \lambda'_t\|_\infty \leq \max_q (\lambda'^{-1}_{t_q} \lambda''_{t_q} \|\mathbf{Q}_{t, :, q}\|_{S^1})$

We notice that  $\lambda'_t = \sigma'(\mathbf{W}\mathbf{s}_{t-1} + \mathbf{U}\mathbf{x}_t)$  and  $\mathbf{s}_t = \sigma(\mathbf{W}\mathbf{s}_{t-1} + \mathbf{U}\mathbf{x}_t)$ , so according to the chain rule of differential calculus, we have the following equation.

$$\frac{\partial \mathbf{s}_t}{\partial \mathbf{U}} = \frac{\partial \lambda'_t}{\partial \mathbf{U}} \rightarrow \frac{\partial \lambda'_t}{\partial \mathbf{U}} = \mathbf{Q}_t \tilde{\Lambda}'^{-1}_t \tilde{\Lambda}''_t \tag{19}$$

So we can derive the upper bound of p-th element of  $\frac{d}{dt} \lambda'_t$  and this is consistent with **Lemma 3**. #

$$[\frac{d}{dt} \lambda_t]_p = tr[\frac{\partial \lambda'_t}{\partial \mathbf{W}} \Delta\mathbf{W}] \leq \|\mathbf{P}_{t-1, :, p} \lambda'^{-1}_p \lambda''_p\|_{S^1} \|\Delta\mathbf{W}\|_{S^\infty} = \lambda'^{-1}_p \lambda''_p \|\mathbf{P}_{t-1, :, p}\|_{S^1} \|\Delta\mathbf{W}\|_{S^\infty} \tag{20}$$

**Lemma 4** For a 3-d tensor  $\mathbf{T}$  and a vector  $\mathbf{v}$ , we have  $\|\mathbf{T}\tilde{\mathbf{v}}\|_{S^1} \leq \|\mathbf{T}_{:, :, *}\|_{S^1} \|\mathbf{v}\|_2$

This lemma can be simply proved using the triangle property of Schatten-1 and Euclidean norm.

$$\|\mathbf{T}\tilde{\mathbf{v}}\|_{S^1} = \|\sum_p \mathbf{v}_p \mathbf{T}_{:, :, p}\|_{S^1} \leq \sum_p \mathbf{v}_p \|\mathbf{T}_{:, :, p}\|_{S^1} \leq \|\mathbf{T}_{:, :, *}\|_{S^1} \|\mathbf{v}\|_2 \# \tag{21}$$

Integrate formula 13,14,15,17 and 18, we can obtain the upper bound of Lipschitz constant in Schatten- $\infty$  norm for input matrix  $\mathbf{U}$ .

$$L_{\mathbf{U}} = \frac{1}{2} \max_q \|\mathbf{Q}_t \tilde{\mathbf{V}}_q\|_{S^1}^2 + \max_q \|\mathbf{V}_q\|_2 \|\mathbf{x}_t\|_2 \max_p (\lambda_p'^{-1} \lambda_p'' \|\mathbf{Q}_{t-1:,*,p}\|_{S^1}) \\ + \max_q \|\mathbf{V}_q\|_2 \|\mathbf{W}\|_{S^\infty} \|\|\mathbf{Q}_{t-1:,*,*}\|_{S^1}\|_2 \max_p (\lambda_p'^{-1} \lambda_p'' \|\mathbf{Q}_{t-1:,*,p}\|_{S^1}) \quad (22)$$

## B.2 Lipschitz Constant of the Recurrent Matrix

Like the input matrix, we need to upper bound  $\|\mathbf{V}_{\mathbf{s}_t}(\mathbf{W}^k + \Delta \mathbf{W}) - \mathbf{V}_{\mathbf{s}_t}(\mathbf{W}^k)\|_\infty$  and  $\|\mathbf{V}_{\mathbf{s}_t}(\mathbf{W}^k + \Delta \mathbf{W}) - \mathbf{V}_{\mathbf{s}_t}(\mathbf{W}^k) - \langle \mathbf{V} \nabla_{\mathbf{W}} \mathbf{s}_t(\mathbf{W}^k), \Delta \mathbf{W} \rangle\|_\infty$ .

Like formula 13, we can use first order approximation to obtain the upper bound of the first part.

$$\|\mathbf{V}_{\mathbf{s}_t}(\mathbf{W}^k + \Delta \mathbf{W}) - \mathbf{V}_{\mathbf{s}_t}(\mathbf{W}^k)\|_\infty \leq \max_p \|\mathbf{P}_t \tilde{\mathbf{V}}_p\|_{S^1} \|\Delta \mathbf{W}\|_{S^\infty} \quad (23)$$

For the second part, like analysis for matrix  $\mathbf{U}$ , we have the following formulas.

$$[\mathbf{V}_{\mathbf{s}_t}(\mathbf{W}^k + \Delta \mathbf{W}) - \mathbf{V}_{\mathbf{s}_t}(\mathbf{W}^k) - \langle \mathbf{V} \nabla_{\mathbf{W}} \mathbf{s}_t(\mathbf{W}^k), \Delta \mathbf{W} \rangle]_p \\ = \int_0^1 \text{tr}([\mathbf{V}_{p,:} \nabla_{\mathbf{W}} \mathbf{s}_t(\mathbf{W}^k + t\Delta \mathbf{W}) - \mathbf{V}_{p,:} \nabla_{\mathbf{W}} \mathbf{s}_t(\mathbf{W}^k)] \Delta \mathbf{W}) dt \\ \leq \|\Delta \mathbf{W}\|_{S^\infty} \int_0^1 \|\mathbf{V}_{p,:} \nabla_{\mathbf{W}} \mathbf{s}_t(\mathbf{W}^k + t\Delta \mathbf{W}) - \mathbf{V}_{p,:} \nabla_{\mathbf{W}} \mathbf{s}_t(\mathbf{W}^k)\|_{S^1} dt \quad (24)$$

$$\int_0^1 \|\mathbf{V}_{p,:} \nabla_{\mathbf{W}} \mathbf{s}_t(\mathbf{W}^k + t\Delta \mathbf{W}) - \mathbf{V}_{p,:} \nabla_{\mathbf{W}} \mathbf{s}_t(\mathbf{W}^k)\|_{S^1} dt \\ \simeq \int_0^1 \frac{d}{dt} \|\mathbf{V}_{p,:} \nabla_{\mathbf{W}} \mathbf{s}_t(\mathbf{W}^k + t\Delta \mathbf{W})\|_{S^1} |_{t=0} t dt \\ = \frac{1}{2} \frac{d}{dt} \|\mathbf{V}_{p,:} \nabla_{\mathbf{W}} \mathbf{s}_t(\mathbf{W}^k + t\Delta \mathbf{W})\|_{S^1} |_{t=0} \quad (25)$$

Unlike scenarios for  $\mathbf{U}$ ,  $\frac{d}{dt} \|\mathbf{V}_{p,:} \nabla_{\mathbf{W}} \mathbf{s}_t(\mathbf{W}^k + t\Delta \mathbf{W})\|_{S^1} |_{t=0}$  is complicated and contains four parts.

$$\nabla_{\mathbf{W}}(\mathbf{V}_p \mathbf{s}_t) = (\tilde{\mathbf{S}}_{t-1} + \mathbf{P}_{t-1} \tilde{\mathbf{W}}^T) \tilde{\Lambda}'_t \tilde{\mathbf{V}}_p^T = \Lambda'_t \mathbf{V}_p^T \mathbf{s}_{t-1} + \mathbf{P}_{t-1} \tilde{\mathbf{W}}^T \tilde{\Lambda}'_t \tilde{\mathbf{V}}_p^T \quad (26)$$

$$\frac{d}{dt} \nabla_{\mathbf{W}}(\mathbf{V}_p \mathbf{s}_t) = \frac{d}{dt} \Lambda'_t \mathbf{V}_p^T \mathbf{s}_{t-1} + \Lambda'_t \mathbf{V}_p^T \frac{d}{dt} \mathbf{s}_{t-1} + \mathbf{P}_{t-1} \frac{d}{dt} \tilde{\mathbf{W}}^T \tilde{\Lambda}'_t \tilde{\mathbf{V}}_p^T + \mathbf{P}_{t-1} \tilde{\mathbf{W}}^T \frac{d}{dt} \tilde{\Lambda}'_t \tilde{\mathbf{V}}_p^T \quad (27)$$

Similarly, we have to find the upper bound of each part using  $\|\Delta \mathbf{W}\|_{S^\infty}$ .

$$\|\frac{d}{dt} \Lambda'_t \mathbf{V}_p^T \mathbf{s}_{t-1}\|_{S^1} = \|(\frac{d}{dt} \lambda'_t \odot \mathbf{V}_p^T) \mathbf{s}_{t-1}\|_{S^1} \leq \|\frac{d}{dt} \lambda'_t \odot \mathbf{V}_p^T\|_2 \|\mathbf{s}_{t-1}\|_2 \\ \leq \|\frac{d}{dt} \lambda'_t\|_\infty \|\mathbf{V}_p\|_2 \|\mathbf{s}_{t-1}\|_2 = \max_q (\lambda_q'^{-1} \lambda_q'' \|\mathbf{P}_{t-1:,*,q}\|_{S^1}) \|\Delta \mathbf{W}\|_{S^\infty} \|\mathbf{V}_p\|_2 \|\mathbf{s}_{t-1}\|_2 \quad (28)$$

$$\|\Lambda'_t \mathbf{V}_{p,:} (\frac{d}{dt} \mathbf{s}_{t-1}^T)\|_{S^1} \leq \|\lambda'_t \odot \mathbf{V}_p^T\|_2 \|\frac{d}{dt} \mathbf{s}_{t-1}\|_2 \leq \|\lambda'_t\|_\infty \|\mathbf{V}_p\|_2 \|\text{tr}[\frac{\partial \mathbf{s}_{t-1}}{\partial \mathbf{W}} \Delta \mathbf{W}]\|_2 \\ = \|\lambda'_t\|_\infty \|\mathbf{V}_p\|_2 \|\Delta \mathbf{W}\|_{S^\infty} \|\|\mathbf{P}_{t-1:,*,*}\|_{S^1}\|_2 \quad (29)$$

$$\begin{aligned}
& \|\mathbf{P}_{t-1} \frac{d}{dt} \tilde{\mathbf{W}}^T \tilde{\mathbf{\Lambda}}'_t \tilde{\mathbf{V}}_p^T\|_{S^1} \leq \| \|\mathbf{P}_{t-1, :, *}\|_{S^1} \|_2 \| \frac{d}{dt} \mathbf{W}^T \mathbf{\Lambda}'_t \mathbf{V}_p^T \|_2 \\
& = \| \|\mathbf{P}_{t-1, :, *}\|_{S^1} \|_2 \| \Delta \mathbf{W}^T \mathbf{\Lambda}'_t \mathbf{V}_p^T \| \| \leq \| \|\mathbf{P}_{t-1, :, *}\|_{S^1} \|_2 \| \Delta \mathbf{W} \|_{S^\infty} \| \lambda'_t \odot \mathbf{V}_p^T \|_2 \\
& \leq \| \|\mathbf{P}_{t-1, :, *}\|_{S^1} \|_2 \| \Delta \mathbf{W} \|_{S^\infty} \| \lambda'_t \|_\infty \| \mathbf{V}_p^T \|_2
\end{aligned} \tag{30}$$

$$\begin{aligned}
& \|\mathbf{P}_{t-1} \tilde{\mathbf{W}}^T \frac{d}{dt} \tilde{\mathbf{\Lambda}}'_t \tilde{\mathbf{V}}_p^T\|_{S^1} \leq \| \|\mathbf{P}_{t-1, :, *}\|_{S^1} \|_2 \| \mathbf{W}^T \frac{d}{dt} \mathbf{\Lambda}'_t \mathbf{V}_p^T \|_2 \\
& \leq \| \|\mathbf{P}_{t-1, :, *}\|_{S^1} \|_2 \| \mathbf{W} \|_{S^\infty} \| \frac{d}{dt} \lambda'_t \odot \mathbf{V}_p^T \| \leq \| \|\mathbf{P}_{t-1, :, *}\|_{S^1} \|_2 \| \mathbf{W} \|_{S^\infty} \| \frac{d}{dt} \lambda'_t \|_\infty \| \mathbf{V}_p \|_2 \\
& \leq \| \|\mathbf{P}_{t-1, :, *}\|_{S^1} \|_2 \| \mathbf{W} \|_{S^\infty} \max_q (\lambda_{t_q}'^{-1} \lambda_{t_q}'' \| \mathbf{P}_{t-1, :, q} \|) \| \Delta \mathbf{W} \|_{S^\infty} \| \mathbf{V}_p \|_2
\end{aligned} \tag{31}$$

Combined formulas above, we can obtain the Lipschitz constant in Schatten- $\infty$  norm for recurrent matrix  $\mathbf{W}$ .

$$\begin{aligned}
L_{\mathbf{W}} &= \frac{1}{2} (\max_q \| \mathbf{P}_t \tilde{\mathbf{V}}_q \|_{S^1})^2 + \max_q \| \mathbf{V}_q \|_2 \| \mathbf{s}_{t-1} \|_2 \max_p (\lambda_p'^{-1} \lambda_p'' \| \mathbf{P}_{t-1, :, p} \|_{S^1}) \\
&+ 2 \max_q \| \mathbf{V}_q \|_2 \| \lambda'_t \|_\infty \| \|\mathbf{P}_{t-1, :, *}\|_{S^1} \|_2 + \max_q \| \mathbf{V}_q \|_2 \| \mathbf{W} \|_{S^\infty} \| \|\mathbf{P}_{t-1, :, *}\|_{S^1} \|_2 \max_p (\lambda_p'^{-1} \lambda_p'' \| \mathbf{P}_{t-1, :, p} \|)
\end{aligned} \tag{32}$$

## C References

- [1]. Carlson, David E., Volkan Cevher, and Lawrence Carin. "Stochastic Spectral Descent for Restricted Boltzmann Machines." AISTATS. 2015.
- [2]. Bernstein, Dennis S. *Matrix mathematics: theory, facts, and formulas*. Princeton University Press, 2009.
- [3]. Carlson, David E., et al. "Preconditioned spectral descent for deep learning." Advances in Neural Information Processing Systems. 2015.