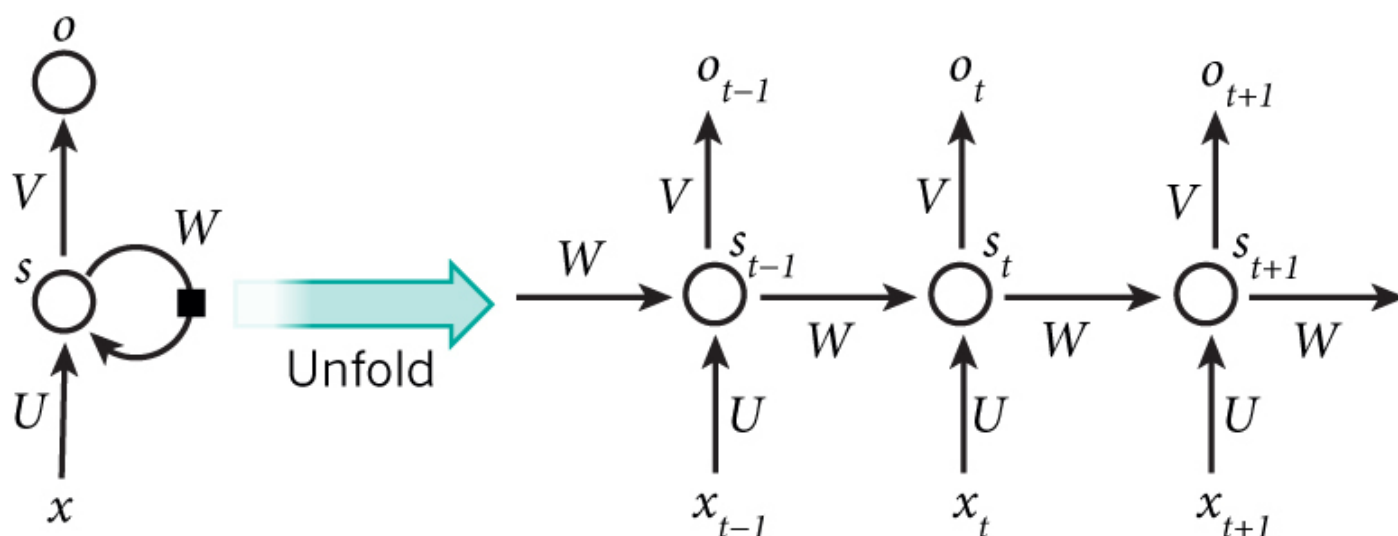


# Summary of Last Month (2.22-3.20)

Liu, Chen. 20th, March

## Basic Knowledge of RNN

---



### 1. Notation

- Input data:  $\mathbf{x} \in \mathbb{R}^N$
- Hidden state:  $\mathbf{s} \in \mathbb{R}^H$
- Output data:  $\mathbf{o} \in \mathbb{R}^K$
- Input connection:  $\mathbf{U} \in \mathbb{R}^{H \times N}$
- Recurrent connection:  $\mathbf{W} \in \mathbb{R}^{H \times H}$
- Output connection:  $\mathbf{V} \in \mathbb{R}^{H \times K}$

### 2. Computation:

- $\mathbf{s}_t = \sigma(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{s}_{t-1})$
- $\mathbf{o}_t = \text{softmax}(\mathbf{V}\mathbf{s}_t)$

### 3. Loss function:

- Usually, we treat a sequence  $\{\mathbf{x}_i\}_{1 \leq i \leq M}$  as a training instance. Given the label of such sequence  $\{\mathbf{y}_i\}_{1 \leq i \leq M}$ , the loss function is defined below:

$$E = -\frac{1}{M} \sum_{t=1}^M \mathbf{y}_t^T \log \mathbf{o}_t$$

## Loss Function of RNN

---

1. Obviously, we have that all the entries in  $\mathbf{y}_t$  sum up to one. We can simplify the expression of loss function of RNN.

$$\begin{aligned}
E &= -\frac{1}{M} \sum_{t=1}^M \sum_{i=1}^K \mathbf{y}_{t,i} \log \mathbf{o}_{t,i} \\
&= -\frac{1}{M} \sum_{t=1}^M \sum_{i=1}^K \mathbf{y}_{t,i} \log \frac{e^{\mathbf{V}_{i,:} \mathbf{s}_t}}{\sum_{j=1}^K e^{\mathbf{V}_{j,:} \mathbf{s}_t}} \\
&= \frac{1}{M} \sum_{t=1}^M (\log(\sum_{j=1}^K e^{\mathbf{V}_{j,:} \mathbf{s}_t}) - \mathbf{y}_t^T \mathbf{V} \mathbf{s}_t)
\end{aligned}$$

2. We define  $\theta = \{\mathbf{U}, \mathbf{W}, \mathbf{V}\}$ ,  $f(\theta) = \sum_{t=1}^M \log(\sum_{j=1}^K e^{\mathbf{V}_{j,:} \mathbf{s}_t})$  and  $g(\theta) = \sum_{t=1}^M \mathbf{y}_t^T \mathbf{V} \mathbf{s}_t$ , so

$$E = \frac{1}{M} (f(\theta) - g(\theta))$$

## Loss Bounded with Respect to V

---

1. For  $g(\theta)$  is linear with respect to  $\mathbf{V}$ , so we have:

$$g(\mathbf{V}^k + \Delta \mathbf{V}) = g(\mathbf{V}^k) + \langle \nabla_{\mathbf{V}} g(\mathbf{V}^k), \Delta \mathbf{V} \rangle$$

2.  $f(\theta)$  is sum of several exponentiations. According to [Theorem 1 in "Stochastic Spectral Descent for Restricted Boltzmann Machine"](#). Let  $\text{lse}_{\mathbf{w}}(\mathbf{u}) = \log \sum_{j=1}^K w_j \exp(u_j)$ , we have:

$$f(\mathbf{V}) = \text{lse}_1(\mathbf{V} \mathbf{s}_t)$$

Thus we can draw an upper bound with respect to  $\mathbf{V} \mathbf{s}_t$ :

$$\begin{aligned}
f(\mathbf{U}^k, \mathbf{W}^k, \mathbf{V}^k + \Delta \mathbf{V}) &= \text{lse}_1(\mathbf{V}^k \mathbf{s}_t + \Delta \mathbf{V} \mathbf{s}_t) \\
&\leq \text{lse}_1(\mathbf{V}^k \mathbf{s}_t) + \langle \nabla \text{lse}_1(\mathbf{V}^k \mathbf{s}_t), \Delta \mathbf{V} \mathbf{s}_t \rangle + \frac{1}{2} \|\Delta \mathbf{V} \mathbf{s}_t\|_{\infty}^2 \\
&= f(\theta^k) + \langle \nabla_{\mathbf{V} \mathbf{s}_t} f(\theta^k), \Delta \mathbf{V} \mathbf{s}_t \rangle + \frac{1}{2} \|\Delta \mathbf{V} \mathbf{s}_t\|_{\infty}^2 \\
&\leq f(\theta^k) + \langle \nabla_{\mathbf{V}} f(\theta^k), \Delta \mathbf{V} \rangle + \frac{1}{2} \|\Delta \mathbf{V}\|_{S^{\infty}}^2 \|\mathbf{s}_t\|_2^2
\end{aligned}$$

3.  $E(\mathbf{U}^k, \mathbf{W}^k, \mathbf{V}^k + \Delta \mathbf{V}) \leq E(\theta^k) + \langle \nabla_{\text{vec}(\mathbf{V})} E(\theta^k), \text{vec}(\Delta \mathbf{V}) \rangle + \frac{1}{2} \|\Delta \mathbf{V}\|_{S^{\infty}}^2 \|\mathbf{s}_t\|_2^2$

## Loss Bound with Respect to W

---

1. Let  $\Phi = \{\mathbf{W}, \mathbf{U}\}$ , we can regard  $\mathbf{s}_t$  as a function of  $\Phi$ :  $\mathbf{s}_t(\Phi)$ . We have an upper bound of  $E(\Phi + \Delta \Phi)$ :  
(Proof of this equation are attached in the handcraft appendix)

$$E(\Phi + \Delta\Phi) \leq E(\Phi) + \langle \nabla_{\Phi} E(\Phi), \Delta\Phi \rangle + \frac{1}{M} \sum_{t=1}^M \left( \frac{1}{2} \|\mathbf{V}\mathbf{s}_t(\Phi + \Delta\Phi) - \mathbf{V}\mathbf{s}_t(\Phi)\|_{\infty}^2 + 2\|\mathbf{V}\mathbf{s}_t(\Phi + \Delta\Phi) - \mathbf{V}\mathbf{s}_t(\Phi) - \langle \mathbf{V}\nabla_{\mathbf{s}_t(\Phi)}, \Delta\Phi \rangle\|_{\infty} \right)$$

2. The upper bound of  $\|\mathbf{V}\mathbf{s}_t(\mathbf{W}^k + \Delta\mathbf{W}) - \mathbf{V}\mathbf{s}_t(\mathbf{W}^k)\|_{\infty}^2$  to update  $\mathbf{W}$ .

1. Last week, I made a mistake.  $\frac{\partial[\mathbf{s}_t(\mathbf{W})]_p}{\partial \mathbf{W}_{ij}} \neq 0$  even if  $p \neq i$  when  $t > 2$ . So the linking equation is:

$$\frac{[\partial \mathbf{V}\mathbf{s}_t(\mathbf{W})]_p}{\partial \mathbf{W}_{ij}} = \sum_{k=1}^H \mathbf{V}_{pk} \frac{\partial[\mathbf{s}_t(\mathbf{W})]_k}{\partial \mathbf{W}_{ij}}$$

$$\frac{\partial[\mathbf{s}_t(\mathbf{W})]_p}{\partial \mathbf{W}_{ij}} = \sigma'(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{s}_{t-1})_p ([\mathbf{s}_{t-1}]_j \delta_{ip} + \sum_{k=1}^H \frac{\partial[\mathbf{s}_{t-1}(\mathbf{W})]_k}{\partial \mathbf{W}_{ij}} \mathbf{W}_{pk})$$

where  $\delta_{ij} = 1$  iff  $i = j$ .

2. Let  $P_t$  be a matrix of size  $H \times H^2$  where  $P_t(k, iH + j) = \frac{\partial[\mathbf{s}_t(\mathbf{W})]_k}{\partial \mathbf{W}_{ij}}$ ,  $\bar{\mathbf{S}}_t$  be a sparse matrix of the same size where  $\bar{\mathbf{S}}_t(i, (i-1)H + j) = \mathbf{s}_{t,j}$  and  $\Lambda_t$  be a diagonal matrix whose diagonal is vector  $\sigma'(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{s}_{t-1})$ . We have:

$$\frac{[\partial \mathbf{V}\mathbf{s}_t(\mathbf{W})]_p}{\partial \mathbf{W}} = Mtr(\mathbf{V}_{p,:} \mathbf{P}_t)$$

$$\mathbf{P}_t = \Lambda_t (\bar{\mathbf{S}}_{t-1} + \mathbf{W}\mathbf{P}_{t-1})$$

$Mtr(\mathbf{s})$  is to turn a  $H^2$ -dim vector into a  $H \times H$  matrix.

3. Now, we can derive the upper bound of  $\|\mathbf{V}\mathbf{s}_t(\mathbf{W}^k + \Delta\mathbf{W}) - \mathbf{V}\mathbf{s}_t(\mathbf{W}^k)\|_{\infty}^2$

$$\begin{aligned} [\mathbf{V}\mathbf{s}_t(\mathbf{W}^k + \Delta\mathbf{W}) - \mathbf{V}\mathbf{s}_t(\mathbf{W}^k)]_p &= \int_0^1 tr \left[ \frac{\partial[\mathbf{V}\mathbf{s}_t(\mathbf{W})]_p}{\partial \mathbf{W}} \Big|_{\mathbf{W}=\mathbf{W}^k + t\Delta\mathbf{W}} \Delta\mathbf{W} \right] dt \\ &= \int_0^1 tr [Mtr(\mathbf{V}_{p,:} \Lambda_t (\bar{\mathbf{S}}_{t-1} + \mathbf{W}\mathbf{P}_{t-1})) \Big|_{\mathbf{W}=\mathbf{W}^k + t\Delta\mathbf{W}} \Delta\mathbf{W}] dt \end{aligned}$$

Approximate this formula at  $t = 0$  and note that every diagonal element of  $\Lambda_t$  is smaller than  $\max_x \sigma'(x)$ , so we have:

$$\begin{aligned} [\mathbf{V}\mathbf{s}_t(\mathbf{W}^k + \Delta\mathbf{W}) - \mathbf{V}\mathbf{s}_t(\mathbf{W}^k)]_p &\leq \max_x \sigma'(x) tr [Mtr(\mathbf{V}_{p,:} (\bar{\mathbf{S}}_{t-1} + \mathbf{W}\mathbf{P}_{t-1})) \Big|_{\mathbf{W}=\mathbf{W}^k} \Delta\mathbf{W}^T] \\ &= \max_x \sigma'(x) \{ tr[\mathbf{V}_{p,:}^T \mathbf{s}_{t-1}^T \Delta\mathbf{W}^T] + tr[Mtr(\mathbf{V}_{p,:} \mathbf{W}^k \mathbf{P}_{t-1}) \Delta\mathbf{W}^T] \} \\ &= \max_x \sigma'(x) \{ \mathbf{V}_{p,:} \Delta\mathbf{W} \mathbf{s}_{t-1} + tr[Mtr(\mathbf{V}_{p,:} \mathbf{W}^k \mathbf{P}_{t-1}) \Delta\mathbf{W}] \} \\ &\leq \max_x \sigma'(x) \{ \|\mathbf{V}_{p,:}\|_2 \|\Delta\mathbf{W}\|_{S^{\infty}} \|\mathbf{s}_{t-1}\|_2 + \|Mtr(\mathbf{V}_{p,:} \mathbf{W}^k \mathbf{P}_{t-1})\|_{S^1} \|\Delta\mathbf{W}\|_{S^{\infty}} \} \\ &= \max_x \sigma'(x) \|\Delta\mathbf{W}\|_{S^{\infty}} \{ \|\mathbf{V}_{p,:}\|_2 \|\mathbf{s}_{t-1}\|_2 + \|Mtr(\mathbf{V}_{p,:} \mathbf{W}^k \mathbf{P}_{t-1})\|_{S^1} \} \end{aligned}$$

3. The upper bound of  $\|\mathbf{V}\mathbf{s}_t(\mathbf{W}^k + \Delta\mathbf{W}) - \mathbf{V}\mathbf{s}_t(\mathbf{W}^k) - \langle \mathbf{V}\nabla_{\mathbf{W}\mathbf{s}_t(\mathbf{W}^k)}, \Delta\mathbf{W} \rangle\|_{\infty}$

1. We analyze the  $p$ -th element of this vector.

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

2. Let's focus on the last integration part. We approximate it at point  $t = 0$ :

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

3. We already have:

$$\nabla_{\mathbf{W}}\mathbf{V}_{p,:}\mathbf{s}_t(\mathbf{W}) = \text{Mtr}[\mathbf{V}_{p,:}\Lambda_t(\bar{\mathbf{S}}_{t-1} + \mathbf{W}\mathbf{P}_{t-1})] = \Lambda_t\mathbf{V}_{p,:}^T\mathbf{s}_{t-1}^T + \text{Mtr}(\mathbf{V}_{p,:}\Lambda_t\mathbf{W}\mathbf{P}_{t-1})$$

so( $\lambda_t$  is the diagonal elements of  $\Lambda_t$ ): **We ignore the second-order derivative part, which vanish much faster then the first order part.**

$$\frac{d}{dt}\nabla_{\mathbf{W}}\mathbf{V}_{p,:}\mathbf{s}_t(\mathbf{W}) = \left(\frac{d}{dt}\lambda_t\right) \odot \mathbf{V}_{p,:}^T\mathbf{s}_{t-1}^T + \Lambda_t\mathbf{V}_{p,:}^T\left(\frac{d}{dt}\mathbf{s}_{t-1}^T\right) + \text{Mtr}(\mathbf{V}_{p,:}\frac{d}{dt}\Lambda_t\mathbf{W}\mathbf{P}_{t-1}) + \text{Mtr}(\mathbf{V}_{p,:}\Lambda_t\Delta\mathbf{W}\mathbf{P}_{t-1})$$

1. Upper bound of  $\left(\frac{d}{dt}\lambda_t\right) \odot \mathbf{V}_{p,:}^T\mathbf{s}_{t-1}^T$

- Let  $\lambda_t''$  be a column vector of  $\sigma''(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{s}_t)$ :

$$\frac{d}{dt}\lambda_t = \lambda_t'' \odot [(\bar{\mathbf{S}}_{t-1} + \mathbf{W}\mathbf{P}_{t-1})\text{vec}(\Delta\mathbf{W})] = \lambda_t'' \odot (\Delta\mathbf{W}\mathbf{s}_{t-1} + \mathbf{W}\mathbf{P}_{t-1}\text{vec}(\Delta\mathbf{W}))$$

- Note that  $\mathbf{W}\mathbf{P}_{t-1}\text{vec}(\Delta\mathbf{W})$  is a column vector whose  $p$ -th element is  $\mathbf{W}_{p,:}\mathbf{P}_{t-1}\text{vec}(\Delta\mathbf{W}) = \text{tr}[\text{Mtr}(\mathbf{W}_{p,:}\mathbf{P}_{t-1})\Delta\mathbf{W}] \leq \|\text{Mtr}(\mathbf{W}_{p,:}\mathbf{P}_{t-1})\|_{S^1} \|\Delta\mathbf{W}\|_{S^\infty}$
- As a result:

$$\begin{aligned}
\left\| \left(\frac{d}{dt}\lambda_t\right) \odot \mathbf{V}_{p,:}^T\mathbf{s}_{t-1}^T \right\|_{S^1} &\leq \left\| \left(\frac{d}{dt}\lambda_t\right) \odot \mathbf{V}_{p,:}^T \right\|_{S^\infty} \|\mathbf{s}_{t-1}^T\|_{S^1} = \left\| \left(\frac{d}{dt}\lambda_t\right) \odot \mathbf{V}_{p,:}^T \right\|_2 \|\mathbf{s}_{t-1}^T\|_2 \\
&\leq \|\lambda_t''\|_2 \|\Delta\mathbf{W}\mathbf{s}_{t-1} + \mathbf{W}\mathbf{P}_{t-1}\text{vec}(\Delta\mathbf{W})\|_2 \|\mathbf{V}_{p,:}^T\|_\infty \|\mathbf{s}_{t-1}^T\|_2 \\
&\leq \|\lambda_t''\|_\infty \|\mathbf{V}_{p,:}^T\|_\infty \|\mathbf{s}_{t-1}^T\|_2 (\|\Delta\mathbf{W}\mathbf{s}_{t-1}\|_2 + \|\mathbf{W}\mathbf{P}_{t-1}\text{vec}(\Delta\mathbf{W})\|_2) \\
&\leq \|\lambda_t''\|_\infty \|\mathbf{V}_{p,:}^T\|_\infty \|\mathbf{s}_{t-1}^T\|_2 (\|\Delta\mathbf{W}\|_{S^\infty} \|\mathbf{s}_{t-1}^T\|_2 + \sqrt{H} \max_p(\|\text{Mtr}(\mathbf{W}_{p,:}\mathbf{P}_{t-1})\|_{S^1}) \|\Delta\mathbf{W}\|_{S^\infty}) \\
&= \|\lambda_t''\|_\infty \|\mathbf{V}_{p,:}^T\|_\infty \|\mathbf{s}_{t-1}^T\|_2 \|\Delta\mathbf{W}\|_{S^\infty} (\|\mathbf{s}_{t-1}^T\|_2 + \sqrt{H} \max_p(\|\text{Mtr}(\mathbf{W}_{p,:}\mathbf{P}_{t-1})\|_{S^1}))
\end{aligned}$$

We have taken advantage of the conclusion: 1)  $\|\mathbf{x} \odot \mathbf{y}\|_2 \leq \|\mathbf{x}\|_\infty \|\mathbf{y}\|_2$ . 2)

$\|\mathbf{A}\mathbf{B}\|_{S^r} \leq \|\mathbf{A}\|_{S^q} \|\mathbf{B}\|_{S^p}$  if  $r^{-1} = p^{-1} + q^{-1}$ . 3) For vector (no matter row or column)  $\mathbf{x}$ , we have  $\|\mathbf{x}\|_{S^\infty} = \|\mathbf{x}\|_{S^1} = \|\mathbf{x}\|_2$ .

2. Upper bound of  $\Lambda_t \mathbf{V}_{p,:}^T (\frac{d}{d_t} \mathbf{s}_{t-1}^T)$

- It is easy to find out that  $\frac{d}{d_t} \mathbf{s}_{t-1}^T = \mathbf{P}_{t-1} \text{vec}(\Delta \mathbf{W})$  whose  $p$ -th element is bounded by  $\|Mtr(\mathbf{P}_{t-1,p})\|_{S^1} \|\Delta \mathbf{W}\|_{S^\infty}$
- Same as above:

$$\begin{aligned} \|\Lambda_t \mathbf{V}_{p,:}^T (\frac{d}{d_t} \mathbf{s}_{t-1}^T)\|_{S^1} &\leq \|\Lambda_t \mathbf{V}_{p,:}^T\|_2 \|\mathbf{P}_{t-1} \text{vec}(\Delta \mathbf{W})\|_2 \\ &\leq \|\lambda_t\|_\infty \|\mathbf{V}_{p,:}^T\|_2 \|\Delta \mathbf{W}\|_{S^\infty} \sqrt{H} \max_p (\|Mtr(\mathbf{P}_{t-1,p})\|_{S^1}) \end{aligned}$$

3. Upper bound of  $Mtr(\mathbf{V}_{p,:} \frac{d}{d_t} \Lambda_t \mathbf{W} \mathbf{P}_{t-1})$

- We first estimate the upper bound of each element and then estimate the upper bound of Schatten-1 norm. Let  $\mathbf{M} = Mtr(\mathbf{V}_{p,:} \frac{d}{d_t} \Lambda_t \mathbf{W} \mathbf{P}_{t-1})$ . For example, we can bound  $\mathbf{M}_{ij}$  by:

$$\begin{aligned} \mathbf{M}_{ij} &= (\frac{d}{d_t} \lambda_t^T) \odot \mathbf{V}_{p,:} \mathbf{W} \mathbf{P}_{t-1, :, iH+j} \leq \|(\frac{d}{d_t} \lambda_t^T) \odot \mathbf{V}_{p,:} \mathbf{W}\|_2 \|\mathbf{P}_{t-1, :, iH+j}\|_2 \\ &\leq \|\frac{d}{d_t} \lambda_t^T\|_2 \|\mathbf{V}_{p,:}\|_\infty \|\mathbf{W}\|_{S^\infty} \|\mathbf{P}_{t-1, :, iH+j}\|_2 \\ &\leq \|\lambda_t''\|_\infty \|\Delta \mathbf{W}\|_{S^\infty} (\|\mathbf{s}_{t-1}^T\|_2 + \sqrt{H} \max_p (\|Mtr(\mathbf{W}_{p,:} \mathbf{P}_{t-1})\|_{S^1})) \|\mathbf{V}_{p,:}\|_\infty \|\mathbf{W}\|_{S^\infty} \|\mathbf{P}_{t-1, :, iH+j}\|_2 \end{aligned}$$

- Now we estimate a upper bound of each element of  $\mathbf{M}$ . Note that Schatten-1 norm and element-wise-1 norm are equivalent,  $\frac{1}{H} \|\mathbf{M}\| \leq \|\mathbf{M}\|_{S^\infty} \leq \|\mathbf{M}\|$ , so we can approximately think that a matrix whose elements are all greater than other one has a greater Schatten-1 norm.

$$\|\mathbf{M}\|_{S^1} \leq \|\lambda_t''\|_\infty \|\Delta \mathbf{W}\|_{S^\infty} (\|\mathbf{s}_{t-1}^T\|_2 + \sqrt{H} \max_p (\|Mtr(\mathbf{W}_{p,:} \mathbf{P}_{t-1})\|_{S^1})) \|\mathbf{V}_{p,:}\|_\infty \|\mathbf{W}\|_{S^\infty} \|\bar{\mathbf{P}}_{t-1}\|_{S^1}$$

$$\text{where } \bar{\mathbf{P}}_{t-1}(i,j) = \|\mathbf{P}_{t-1}(:, iH+j)\|_2.$$

4. Upper bound of  $Mtr(\mathbf{V}_{p,:} \Lambda_t \Delta \mathbf{W} \mathbf{P}_{t-1})$

- Very similar to above, first estimate the upper bound of each element.

$$\mathbf{M}_{ij} = \lambda_t^T \odot \mathbf{V}_{p,:} \Delta \mathbf{W} \mathbf{P}_{t-1, :, iH+j} \leq \|\lambda_t\|_\infty \|\mathbf{V}_{p,:}\|_2 \|\Delta \mathbf{W}\|_{S^\infty} \|\mathbf{P}_{t-1, :, iH+j}\|_2$$

- So we can get the upper bound:

$$\|\mathbf{M}\|_{S^1} \leq \|\lambda_t\|_\infty \|\mathbf{V}_{p,:}\|_2 \|\Delta \mathbf{W}\|_{S^\infty} \|\bar{\mathbf{P}}_{t-1}\|_{S^1}$$

## Loss Bound with Respect to U

- The upper bound of  $\|\mathbf{V}_{s_i}(\mathbf{U}^k + \Delta \mathbf{U}) - \mathbf{V}_{s_i}(\mathbf{U}^k)\|_\infty^2$  to update  $\mathbf{U}$ .

- Linking equation:

$$\frac{\partial[\mathbf{V}\mathbf{s}_t(\mathbf{U})]_p}{\partial \mathbf{U}_{ij}} = \sum_{k=1}^H \mathbf{V}_{pk} \frac{\partial[\mathbf{s}_t(\mathbf{U})]_k}{\partial \mathbf{U}_{ij}}$$

$$\frac{\partial[\mathbf{s}_t(\mathbf{U})]_p}{\partial \mathbf{U}_{ij}} = \sigma'(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{s}_{t-1})_p [\mathbf{x}_{t,j} \delta_{ip} + \sum_{k=1}^H \mathbf{W}_{pk} \frac{\partial[\mathbf{s}_{t-1}(\mathbf{U})]_k}{\partial \mathbf{U}_{ij}}]$$

Similarly, we have the matrix form:

$$\frac{\partial[\mathbf{V}\mathbf{s}_t(\mathbf{U})]_p}{\partial \mathbf{U}} = \text{Mtr}(\mathbf{V}_{p,:} \mathbf{Q}_t)$$

$$\mathbf{Q}_t = \mathbf{\Lambda}_t(\bar{\mathbf{X}}_t + \mathbf{W}\mathbf{Q}_{t-1})$$

Similar to updating  $\mathbf{W}$ ,  $\mathbf{Q}_t$  is a  $H \times HN$  matrix where  $\mathbf{Q}_t(p, iN + j) = \frac{\partial[\mathbf{s}_t(\mathbf{U})]_p}{\partial \mathbf{U}_{ij}}$ ,  $\bar{\mathbf{X}}_t$  is a sparse matrix of the same size where  $\bar{\mathbf{X}}_t(i, iN + j) = \mathbf{x}_t(j)$  and  $\mathbf{\Lambda}$  is a diagonal derivation matrix.

2. We can now get the upper bound:

$$[\mathbf{V}\mathbf{s}_t(\mathbf{U}^k + \Delta\mathbf{U}) - \mathbf{V}\mathbf{s}_t(\mathbf{U}^k)]_p = \int_0^1 \text{tr}[\frac{\partial[\mathbf{V}\mathbf{s}_t(\mathbf{U})]_p}{\partial \mathbf{U}}]_{\mathbf{U}=\mathbf{U}^k+t\Delta\mathbf{U}} \Delta\mathbf{U} d_t$$

$$= \int_0^1 \text{tr}[\text{Mtr}(\mathbf{V}_{p,:} \mathbf{\Lambda}_t(\bar{\mathbf{X}}_t + \mathbf{W}\mathbf{Q}_{t-1}))]_{\mathbf{U}=\mathbf{U}^k+t\Delta\mathbf{U}} \Delta\mathbf{U} d_t$$

Approximate this equation at  $t = 0$ :

$$[\mathbf{V}\mathbf{s}_t(\mathbf{U}^k + \Delta\mathbf{U}) - \mathbf{V}\mathbf{s}_t(\mathbf{U}^k)]_p \leq \max_x \sigma'(x) \{ \mathbf{V}_{p,:} \Delta\mathbf{U} \mathbf{x}_t + \text{tr}[\text{Mtr}(\mathbf{V}_{p,:} \mathbf{W}^k \mathbf{Q}_{t-1}) \Delta\mathbf{U}] \}$$

$$\leq \max_x \sigma'(x) \{ \|\mathbf{V}_{p,:}\|_2 \|\Delta\mathbf{U}\|_{S^\infty} \|\mathbf{x}_t\|_2 + \|\text{Mtr}(\mathbf{V}_{p,:} \mathbf{W}^k \mathbf{Q}_{t-1})\|_{S^1} \|\Delta\mathbf{U}\|_{S^\infty} \}$$

$$= \max_x \sigma'(x) \|\Delta\mathbf{U}\|_{S^\infty} \{ \|\mathbf{V}_{p,:}\|_2 \|\mathbf{x}_t\|_2 + \|\text{Mtr}(\mathbf{V}_{p,:} \mathbf{W}^k \mathbf{Q}_{t-1})\|_{S^1} \}$$

2. The upper bound of  $\|\mathbf{V}\mathbf{s}_t(\mathbf{U}^k + \Delta\mathbf{U}) - \mathbf{V}\mathbf{s}_t(\mathbf{U}^k) - \langle \mathbf{V} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k), \Delta\mathbf{U} \rangle\|_\infty$

1. Analyze the  $p$ -th element

$$[\mathbf{V}\mathbf{s}_t(\mathbf{U}^k + \Delta\mathbf{U}) - \mathbf{V}\mathbf{s}_t(\mathbf{U}^k) - \langle \mathbf{V} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k), \Delta\mathbf{U} \rangle]_p$$

$$= \int_0^1 \text{tr}([\mathbf{V}_{p,:} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k + t\Delta\mathbf{U}) - \mathbf{V}_{p,:} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k)] \Delta\mathbf{U}) d_t$$

$$\leq \|\Delta\mathbf{U}\|_{S^\infty} \left\| \int_0^1 \mathbf{V}_{p,:} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k + t\Delta\mathbf{U}) - \mathbf{V}_{p,:} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k) d_t \right\|_{S^1}$$

2. Similarly, focus on the last integration part:

$$\begin{aligned}
& \int_0^1 \mathbf{V}_{p,:} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k + t\Delta\mathbf{U}) - \mathbf{V}_{p,:} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k) \\
& \simeq \int_0^1 \frac{d}{dt} \mathbf{V}_{p,:} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k + t\Delta\mathbf{U})|_{t=0} dt \\
& = \frac{1}{2} \frac{d}{dt} \mathbf{V}_{p,:} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k + t\Delta\mathbf{U})|_{t=0} \\
& = \frac{1}{2} \frac{d}{dt} (\lambda_t \odot \mathbf{V}_{p,:}^T \mathbf{x}_t^T + \text{Mtr}(\lambda_t \odot \mathbf{V}_{p,:} \mathbf{W} \mathbf{Q}_{t-1}))|_{t=0} \\
& \simeq \frac{1}{2} \left( \frac{d}{dt} \lambda_t \odot \mathbf{V}_{p,:}^T \mathbf{x}_t^T + \text{Mtr}\left(\frac{d}{dt} \lambda_t \odot \mathbf{V}_{p,:} \mathbf{W} \mathbf{Q}_{t-1}\right) \right)
\end{aligned}$$

In this situation, we have ignored the second derivation part.

3. Upper bound of  $\left\| \frac{d}{dt} \lambda_t \odot \mathbf{V}_{p,:}^T \mathbf{x}_t^T \right\|_{S^1}$

- Let  $\lambda_t''$  be a column vector of  $\sigma''(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{s}_{t-1})$ :

$$\begin{aligned}
\frac{d}{dt} \lambda_t &= \lambda_t'' \odot [(\bar{\mathbf{X}}_t + \mathbf{W}\mathbf{Q}_{t-1})\text{vec}(\Delta\mathbf{U})] = \lambda_t'' \odot (\Delta\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{Q}_{t-1}\text{vec}(\Delta\mathbf{U})) \\
\left\| \left( \frac{d}{dt} \lambda_t \right) \odot \mathbf{V}_{p,:}^T \mathbf{x}_t^T \right\|_{S^1} &\leq \left\| \left( \frac{d}{dt} \lambda_t \right) \odot \mathbf{V}_{p,:}^T \right\|_{S^\infty} \left\| \mathbf{x}_t^T \right\|_{S^1} = \left\| \left( \frac{d}{dt} \lambda_t \right) \odot \mathbf{V}_{p,:}^T \right\|_2 \left\| \mathbf{x}_t^T \right\|_2 \\
&\leq \left\| \lambda_t'' \odot (\Delta\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{Q}_{t-1}\text{vec}(\Delta\mathbf{U})) \right\|_2 \left\| \mathbf{V}_{p,:}^T \right\|_\infty \left\| \mathbf{x}_t^T \right\|_2 \\
&\leq \left\| \lambda_t'' \right\|_\infty \left\| \mathbf{V}_{p,:}^T \right\|_\infty \left\| \mathbf{x}_t^T \right\|_2 (\left\| \Delta\mathbf{U}\mathbf{x}_t \right\|_2 + \left\| \mathbf{W}\mathbf{Q}_{t-1}\text{vec}(\Delta\mathbf{U}) \right\|_2) \\
&\leq \left\| \lambda_t'' \right\|_\infty \left\| \mathbf{V}_{p,:}^T \right\|_\infty \left\| \mathbf{x}_t^T \right\|_2 (\left\| \Delta\mathbf{U} \right\|_{S^\infty} \left\| \mathbf{x}_t^2 \right\|_2 + \sqrt{H} \max_p (\left\| \text{Mtr}(\mathbf{W}_{p,:} \mathbf{Q}_{t-1}) \right\|_{S^1}) \left\| \Delta\mathbf{U} \right\|_{S^\infty}) \\
&= \left\| \lambda_t'' \right\|_\infty \left\| \mathbf{V}_{p,:}^T \right\|_\infty \left\| \mathbf{x}_t^T \right\|_2 \left\| \Delta\mathbf{U} \right\|_{S^\infty} (\left\| \mathbf{x}_t^2 \right\|_2 + \sqrt{H} \max_p (\left\| \text{Mtr}(\mathbf{W}_{p,:} \mathbf{Q}_{t-1}) \right\|_{S^1}))
\end{aligned}$$

4. Upper bound of  $\text{Mtr}\left(\frac{d}{dt} \lambda_t \odot \mathbf{V}_{p,:} \mathbf{W} \mathbf{Q}_{t-1}\right)$

- Let  $\mathbf{N} = \text{Mtr}\left(\frac{d}{dt} \lambda_t \odot \mathbf{V}_{p,:} \mathbf{W} \mathbf{Q}_{t-1}\right)$ , we have:

$$\begin{aligned}
\mathbf{N}_{ij} &= \frac{d}{dt} \lambda_t \odot \mathbf{V}_{p,:} \mathbf{W} \mathbf{Q}_{t-1, :, iN+j} \leq \left\| \left( \frac{d}{dt} \lambda_t^T \right) \odot \mathbf{V}_{p,:} \mathbf{W} \right\|_2 \left\| \mathbf{Q}_{t-1, :, iN+j} \right\|_2 \\
&\leq \left\| \frac{d}{dt} \lambda_t^T \right\|_2 \left\| \mathbf{V}_{p,:} \right\|_\infty \left\| \mathbf{W} \right\|_{S^\infty} \left\| \mathbf{Q}_{t-1, :, iN+j} \right\|_2 \\
&\leq \left\| \lambda_t'' \right\|_\infty \left\| \Delta\mathbf{U} \right\|_{S^\infty} (\left\| \mathbf{x}_t \right\|_2 + \sqrt{H} \max_p (\left\| \text{Mtr}(\mathbf{W}_{p,:} \mathbf{Q}_{t-1}) \right\|_{S^1}) \left\| \mathbf{V}_{p,:} \right\|_\infty \left\| \mathbf{W} \right\|_{S^\infty} \left\| \mathbf{Q}_{t-1, :, iN+j} \right\|_2
\end{aligned}$$

## Comments

- Now, we have bounded the loss function based on Schatten- $\infty$  norm. I am not sure it is completely right. However, it is obvious to find the following flaws.(Fonts in purple)**

1. I have ignored the second derivative part while bounding

$$\left\| \mathbf{V} \mathbf{s}_t(\mathbf{W}^k + \Delta\mathbf{W}) - \mathbf{V} \mathbf{s}_t(\mathbf{W}^k) - \langle \mathbf{V} \nabla_{\mathbf{W}} \mathbf{s}_t(\mathbf{W}^k), \Delta\mathbf{W} \rangle \right\|_\infty$$
 as well as

$$\left\| \mathbf{V} \mathbf{s}_t(\mathbf{U}^k + \Delta\mathbf{U}) - \mathbf{V} \mathbf{s}_t(\mathbf{U}^k) - \langle \mathbf{V} \nabla_{\mathbf{U}} \mathbf{s}_t(\mathbf{U}^k), \Delta\mathbf{U} \rangle \right\|_\infty,$$
 which is the same as what previous paper

"Preconditioned Spectral Descent for Deep Learning" do. I have written a very simple script (input and

hidden units are both scalar) to stimulate the what RNN works. It is found that in fine-tuned model second-order derivation is indeed smaller than first-order if the activation is sigmoid, if the activation is relu or tanh, the second-order derivation is a bit bigger than the first one(they are of the same magnitude). However, it is only the comparison between  $\frac{\partial Q}{\partial U}$  v.s  $Q(\frac{\partial P}{\partial U}$  v.s.  $P$ ). **At present, I can not prove  $\|Mtr(\mathbf{V}_{p,:}\mathbf{\Lambda}_t\mathbf{W}\frac{d\mathbf{P}_{t-1}}{dt})\|_{S^1}$  is significantly smaller than (but it can't not much bigger than the first-order part) other parts of the derivative. The consequence will be catastrophic if it is not true, because computing second-order derivative is very expensive.**

2. It seems that matrix operation is not enough to solve this optimization problem perfectly. As we see,  $\frac{\partial s_t}{\partial \mathbf{W}}$  (for  $\mathbf{U}$  is similar) is the core element of linking equation which will be used repeatedly for optimization. **2D matrix is not powerfull enough to address this question, I think 3D tensor is better.** However, I do not have enough mathematical tools to solve 3D operation, so I project the last two dimension into one, resulting  $H \times H^2$  matrix. It is expedient and causes much problems, such as repeatedly reshaping matrix. Some 'ugly' parts like  $\sqrt{H} \max_p(\|Mtr(\mathbf{W}_{p,:}\mathbf{P}_{t-1})\|_{S^1})$  arises from this. **I guess there exists a tighter and nicer upper bound if we use 3D tensor instead of 2D matrix.**
3. I assume that if matrix  $\mathbf{A}$ 's every element is bigger than  $\mathbf{B}$ 's corresponding element, then  $\|\mathbf{A}\|_{S^1} \leq \|\mathbf{B}\|_{S^1}$ . It is based on the fact that Schatten-1 norm and element-wise 1-norm are equivalent. ( $\frac{1}{\sqrt{mn}} \|\mathbf{X}\|_{\infty} \leq \|\mathbf{X}\|_{S^1} \leq \|\mathbf{X}\|_{\infty}$  for matrix  $\mathbf{X}$  of size  $n \times m$ ) This relation is strict if we use Schatten-2 norm, but S-2 norm is very expensive to compute.(It can be approximated by **super power** method) So, we approximate it by using S-1 norm. (This problem is caused also by using 2D matrix instead of 3D tensor.)
4.  $\mathbf{U}$  and  $\mathbf{W}$  are bounded by Schatten- $\infty$  norm, while  $\mathbf{V}$  are bounded by element-wise  $\infty$ -norm. I will change it to Schatten- $\infty$  norm later.



$$E = \frac{1}{M} \sum_{t=1}^M \left( \log \sum_{j=1}^K e^{V_{j,:}^T S_t} \right) - y_t^T V S_t.$$

$$f_{-t}(\theta) = \log \sum_{j=1}^K e^{V_{j,:}^T S_t} \quad g_{-t}(\theta) = -y_t^T V S_t$$

$$\text{Let } r = V S_t(\emptyset) - V S_t(\theta) - \langle \nabla_{\theta} V S_t(\theta), \emptyset - \theta \rangle$$

$$\text{Lemma 1: } f_{-t}(\emptyset) \leq f_{-t}(\theta) + \langle \nabla_{\theta} f_{-t}(\theta), \emptyset - \theta \rangle + \frac{1}{2} \|V S_t(\emptyset) - V S_t(\theta)\|_{\infty}^2 + r$$

We already know, for function  $l_{se}(x) = \log \sum_{j=1}^J e^{x_j}$ , we have

$$l_{se}(x_2) \leq l_{se}(x_1) + \langle \nabla_{x_1} l_{se}(x_1), x_2 - x_1 \rangle + \frac{1}{2} \|x_2 - x_1\|_{\infty}^2$$

$$\text{so } f_{-t}(\emptyset) \leq f_{-t}(\theta) + \langle \nabla_{V S_t} f_{-t}(\theta), V S_t(\emptyset) - V S_t(\theta) \rangle + \frac{1}{2} \|V S_t(\emptyset) - V S_t(\theta)\|_{\infty}^2$$

$$\begin{aligned} \text{where } \langle \nabla_{V S_t} f_{-t}(\theta), V S_t(\emptyset) - V S_t(\theta) \rangle &= \nabla_{V S_t} f_{-t}(\theta) [r + \langle \nabla_{\theta} V S_t(\theta), \emptyset - \theta \rangle] \\ &= r \cdot \nabla_{V S_t(\theta)} f_{-t}(\theta) + \langle \nabla_{V S_t(\theta)} f_{-t}(\theta), \nabla_{\theta} V S_t(\theta), \emptyset - \theta \rangle \\ &= r \cdot \nabla_{V S_t(\theta)} f_{-t}(\theta) + \langle \nabla_{\theta} f_{-t}(\theta), \emptyset - \theta \rangle \end{aligned}$$

It can be proved that  $\nabla_{V S_t(\theta)} f_{-t}(\theta) \leq 1$  (softmax)

$$\text{so } f_{-t}(\emptyset) \leq f_{-t}(\theta) + \langle \nabla_{\theta} f_{-t}(\theta), \emptyset - \theta \rangle + \frac{1}{2} \|V S_t(\emptyset) - V S_t(\theta)\|_{\infty}^2 + \max_{\|r\|_{\infty}} r$$

$$\text{Lemma 2: } g_{-t}(\emptyset) \leq g_{-t}(\theta) + \langle \nabla_{\theta} g_{-t}(\theta), \emptyset - \theta \rangle + \max_{\|r\|_{\infty}} r$$

$g_{-t}(\theta)$  is linear to  $V S_t(\theta)$ , so.

$$g_{-t}(\emptyset) = g_{-t}(\theta) + \langle \nabla_{V S_t(\theta)}, g_{-t}(\theta), V S_t(\emptyset) - V S_t(\theta) \rangle$$

$$\begin{aligned} \text{similarly, } \langle \nabla_{V S_t(\theta)}, g_{-t}(\theta), V S_t(\emptyset) - V S_t(\theta) \rangle &= \langle \nabla_{\theta} g_{-t}(\theta), \emptyset - \theta \rangle + \max_{\|r\|_{\infty}} r \\ &\leq \langle \nabla_{\theta} g_{-t}(\theta), \emptyset - \theta \rangle + \max r \end{aligned}$$

$$\text{so } g_{-t}(\emptyset) \leq g_{-t}(\theta) + \langle \nabla_{\theta} g_{-t}(\theta), \emptyset - \theta \rangle + \max_{\|r\|_{\infty}} r$$