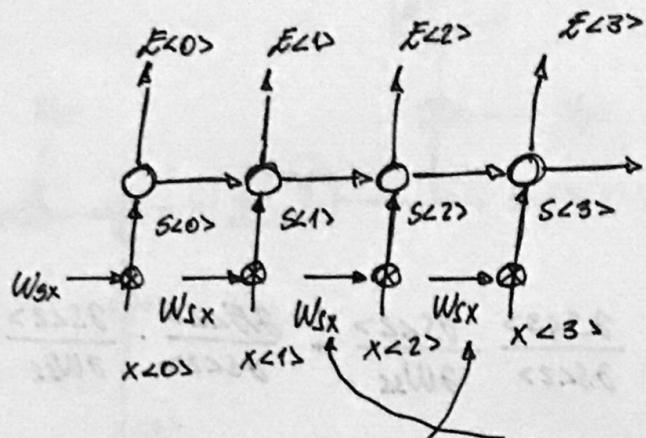
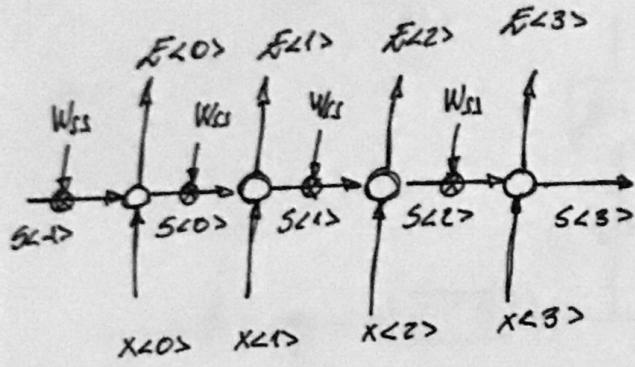


# RNN Network Backprop through time



$$\begin{aligned}
 \frac{\partial E}{\partial W_{sx}} &= \frac{\partial E^{<3>}}{\partial s^{<3>}} \cdot \frac{\partial s^{<3>}}{\partial W_{sx}} + \left[ \frac{\partial E^{<3>}}{\partial s^{<3>}} \cdot \frac{\partial s^{<3>}}{\partial s^{<2>}} \cdot \frac{\partial s^{<2>}}{\partial W_{sx}} + \frac{\partial E^{<2>}}{\partial s^{<2>}} \cdot \frac{\partial s^{<2>}}{\partial W_{sx}} \right. \\
 &\quad + \left[ \frac{\partial E^{<3>}}{\partial s^{<3>}} \cdot \frac{\partial s^{<3>}}{\partial s^{<2>}} \cdot \frac{\partial s^{<2>}}{\partial s^{<1>}} \cdot \frac{\partial s^{<1>}}{\partial W_{sx}} + \frac{\partial E^{<2>}}{\partial s^{<2>}} \cdot \frac{\partial s^{<2>}}{\partial s^{<1>}} \cdot \frac{\partial s^{<1>}}{\partial W_{sx}} + \right. \\
 &\quad \left. + \frac{\partial E^{<1>}}{\partial s^{<1>}} \cdot \frac{\partial s^{<1>}}{\partial W_{sx}} \right] + \left[ \frac{\partial E^{<3>}}{\partial s^{<3>}} \cdot \frac{\partial s^{<3>}}{\partial s^{<2>}} \cdot \frac{\partial s^{<2>}}{\partial s^{<1>}} \cdot \frac{\partial s^{<1>}}{\partial s^{<0>}} \cdot \frac{\partial s^{<0>}}{\partial W_{sx}} + \right. \\
 &\quad \left. + \frac{\partial E^{<2>}}{\partial s^{<2>}} \cdot \frac{\partial s^{<2>}}{\partial s^{<1>}} \cdot \frac{\partial s^{<1>}}{\partial s^{<0>}} \cdot \frac{\partial s^{<0>}}{\partial W_{sx}} + \frac{\partial E^{<1>}}{\partial s^{<1>}} \cdot \frac{\partial s^{<1>}}{\partial s^{<0>}} \cdot \frac{\partial s^{<0>}}{\partial W_{sx}} + \right. \\
 &\quad \left. + \frac{\partial E^{<0>}}{\partial s^{<0>}} \cdot \frac{\partial s^{<0>}}{\partial W_{sx}} \right] = \sum_{K=3}^3 \frac{\partial E^{<K>}}{\partial s^{<K>}} \cdot \prod_{j=3+1}^K \left( \frac{\partial s^{<j>}}{\partial s^{<j-1>}} \right) \cdot \frac{\partial s^{<3>}}{\partial W_{sx}} \\
 &\quad + \sum_{K=2}^3 \frac{\partial E^{<K>}}{\partial s^{<K>}} \cdot \prod_{j=3+1}^K \left( \frac{\partial s^{<j>}}{\partial s^{<j-1>}} \right) \cdot \frac{\partial s^{<2>}}{\partial W_{sx}} + \dots \rightarrow \\
 \frac{\partial E}{\partial W_{sx}} &= \sum_{l=0}^{T_x} \sum_{k=c}^{T_y} \frac{\partial E^{<k>}}{\partial s^{<k>}} \prod_{j=l+1}^K \left( \frac{\partial s^{<j>}}{\partial s^{<j-1>}} \right) \cdot \frac{\partial s^{<l>}}{\partial W_{sx}}
 \end{aligned}$$



$$\begin{aligned} \frac{\partial E}{\partial w_{ss}} &= \left[ \frac{\partial f_{l,3}}{\partial s_{l,3}} \cdot \frac{\partial s_{l,3}}{\partial w_{ss}} \right] + \left[ \frac{\partial f_{l,3}}{\partial s_{l,3}} \cdot \frac{\partial s_{l,3}}{\partial s_{l,2}} \cdot \frac{\partial s_{l,2}}{\partial w_{ss}} + \frac{\partial f_{l,2}}{\partial s_{l,2}} \cdot \frac{\partial s_{l,2}}{\partial w_{ss}} \right] + \\ &+ \left[ \frac{\partial f_{l,3}}{\partial s_{l,3}} \cdot \frac{\partial s_{l,3}}{\partial s_{l,2}} \cdot \frac{\partial s_{l,2}}{\partial s_{l,1}} \cdot \frac{\partial s_{l,1}}{\partial w_{ss}} + \frac{\partial f_{l,2}}{\partial s_{l,2}} \cdot \frac{\partial s_{l,2}}{\partial s_{l,1}} \cdot \frac{\partial s_{l,1}}{\partial w_{ss}} \right. + \\ &\quad \left. + \frac{\partial f_{l,1}}{\partial s_{l,1}} \cdot \frac{\partial s_{l,1}}{\partial w_{ss}} \right] + \left[ \frac{\partial f_{l,0}}{\partial s_{l,0}} \cdot \frac{\partial s_{l,0}}{\partial w_{ss}} \right] + \frac{\partial f_{l,3}}{\partial s_{l,3}} \cdot \frac{\partial s_{l,3}}{\partial s_{l,2}} \cdot \frac{\partial s_{l,2}}{\partial s_{l,1}} \cdot \frac{\partial s_{l,1}}{\partial w_{ss}} \\ &\quad \left. + \frac{\partial f_{l,2}}{\partial s_{l,2}} \cdot \frac{\partial s_{l,2}}{\partial s_{l,1}} \cdot \frac{\partial s_{l,1}}{\partial s_{l,0}} \cdot \frac{\partial s_{l,0}}{\partial w_{ss}} + \frac{\partial f_{l,1}}{\partial s_{l,1}} \cdot \frac{\partial s_{l,1}}{\partial s_{l,0}} \cdot \frac{\partial s_{l,0}}{\partial w_{ss}} \right] = \end{aligned}$$

$$\frac{\partial E}{\partial w_{ss}} = \sum_{l=0}^{T_x} \sum_{k=l}^{T_y} \frac{\partial E_k}{\partial s_{l,k}} \prod_{j=l+1}^K \left( \frac{\partial s_{l,j}}{\partial s_{l,j-1}} \right) \frac{\partial s_{l,0}}{\partial w_{ss}}$$

$$\frac{\partial E}{\partial x_{l,0}} = \sum_{k=l}^{T_y} \frac{\partial E_k}{\partial s_{l,k}} \prod_{j=l+1}^K \left( \frac{\partial s_{l,j}}{\partial s_{l,j-1}} \right) \frac{\partial s_{l,0}}{\partial x_{l,0}}$$

\* These consecutive multiplications may cause vanishing gradient or exploding gradient:

- Exploding gradient  $\rightarrow \frac{\partial E}{\partial \theta} \rightarrow \frac{\partial E}{\partial \theta} > \text{Threshold} \rightarrow \frac{\partial E}{\partial \theta} = \frac{\text{Thres}}{\frac{\partial E}{\partial \theta}} \cdot \frac{\partial E}{\partial \theta}$

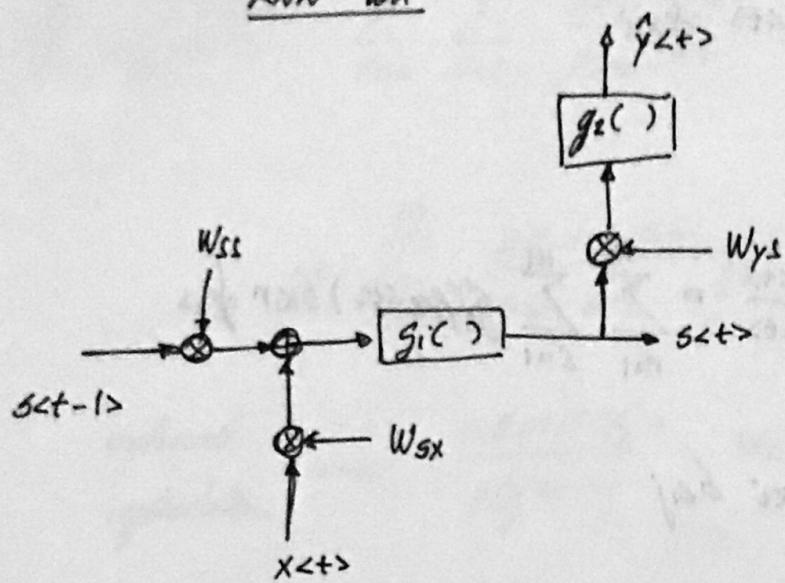
- Vanishing gradient  $\rightarrow \text{LSTM, GRU}$

$$\frac{\partial s_{l,0}}{\partial x_{l,0}} = \tanh'(z_{l,0}) \cdot \text{Weights}_l$$

$$0.x \leftarrow \begin{cases} z = -1 & \\ z = 1 & \end{cases} \leftarrow$$

$\begin{cases} \rightarrow < 1; 0.x^T \rightarrow \text{Vanishing} \\ \rightarrow > 1; 1.x^T \rightarrow \text{Exploding} \end{cases}$

### RNN Cell



\* Forward  $\rightarrow$

$$s<+> = g_1(z<+>)$$

$$z<+> = W_{ss} s<t-1> + W_{sx} x<+> + b_s$$

$$\hat{y}<+> = g_2(p<+>)$$

$$p<+> = W_{ys} s<t> + b_y$$

$$\text{Dimensions} \rightarrow x<+> \in \mathbb{R}^{(n_x, m)}; s<+> \in \mathbb{R}^{(n_h, m)}; y<+> \in \mathbb{R}^{(n_y, m)}$$

$$W_{ss} \in \mathbb{R}^{(n_h, n_h)}; W_{sx} \in \mathbb{R}^{(n_h, n_x)}; W_{ys} \in \mathbb{R}^{(n_y, n_h)}$$

$$\frac{\partial z_{k\ell} <+>}{\partial z_{ij} <+>} = g'_1(z_{k\ell} <+>) \delta_{k,i} \delta_{\ell,j}; \quad \frac{\partial \hat{y}_{k\ell} <+>}{\partial p_{ij} <+>} = g'_2(s_{k\ell} <+>) \delta_{k,i} \delta_{\ell,j}$$

$$\frac{\partial z_{k\ell} <+>}{\partial W_{ss,ij}} = \delta_{k,i} s_{j\ell} <+>; \quad \frac{\partial z_{k\ell} <+>}{\partial W_{sx,ij}} = \delta_{k,i} x_{j\ell} <+>; \quad \frac{\partial z_{k\ell} <+>}{\partial s_{ij} <t-1>} = W_{ss,k\ell} \delta_{\ell,j}$$

$$\frac{\partial z_{k\ell} <+>}{\partial x_{ij} <t>} = W_{sx,ki} \delta_{\ell,j}; \quad \frac{\partial \hat{y}_{k\ell} <+>}{\partial W_{ys,ij}} = \delta_{k,i} s_{j\ell} <+>; \quad \frac{\partial \hat{y}_{k\ell} <+>}{\partial s_{ij} <t>} = W_{ys,ki} \delta_{\ell,j}$$

$$\frac{\partial s_{k\ell} <+>}{\partial x_{ij} <t>} = \sum_{r=1}^{n_h} \sum_{s=1}^m \frac{\partial z_{k\ell} <+>}{\partial z_{rs} <+>} \frac{\partial z_{rs} <+>}{\partial x_{ij} <t>} = \sum_{r=1}^{n_h} \sum_{s=1}^m g'_1(z_{k\ell} <+>) \delta_{k,r} \delta_{\ell,s}$$

$$\cdot W_{sx,ki} \delta_{sj} = g'_1(z_{k\ell} <+>) \cdot W_{sx,ki} \delta_{sj}; \quad \frac{\partial z_{k\ell} <+>}{\partial s_{ij} <t-1>} = g'_1(z_{k\ell} <+>) \cdot W_{ss,ki} \delta_{\ell,j}$$

$$\frac{\partial s_{k\ell} <+>}{\partial W_{ss,ij}} = \sum_{r=1}^{n_h} \sum_{s=1}^m \frac{\partial z_{k\ell} <+>}{\partial z_{rs} <+>} \cdot \frac{\partial z_{rs} <+>}{\partial W_{ss,ij}} = \sum_{r=1}^{n_h} \sum_{s=1}^m g'_1(z_{k\ell} <+>) \delta_{k,r} \delta_{\ell,s} \delta_{ri} \delta_{sj}$$

$$= g'_1(z_{k\ell} <+>) \delta_{j\ell} \delta_{ki} \rightarrow s_{k\ell} <t-1>$$

$$\frac{\partial \hat{y}_{RK}^{re \leftarrow t \rightarrow}}{\partial \cancel{w_{k,i,j}}} = g'_1(z_{RK}^{re \leftarrow t \rightarrow}) \cdot \cancel{x_{je \leftarrow t \rightarrow}} \delta_{K,i}$$

W<sub>k,i,j</sub>

$$\frac{\partial \hat{y}_{RK}^{re \leftarrow t \rightarrow}}{\partial s_{ij}^{re \leftarrow t \rightarrow}} = \sum_{r=1}^{n_r} \sum_{s=1}^m \frac{\partial \hat{y}_{RK}^{re \leftarrow t \rightarrow}}{\partial p_{rs}^{re \leftarrow t \rightarrow}} \cdot \frac{\partial p_{rs}^{re \leftarrow t \rightarrow}}{\partial s_{ij}^{re \leftarrow t \rightarrow}} = \sum_{r=1}^{n_r} \sum_{s=1}^m g'_2(p_{rs}^{re \leftarrow t \rightarrow}) \delta_{K,r} \delta_{e,s}$$

$$W_{p_{rs}^{re \leftarrow t \rightarrow}} \delta_{s,j} = g'_2(p_{rs}^{re \leftarrow t \rightarrow}) \cdot W_{s_{ij}^{re \leftarrow t \rightarrow}} \delta_{e,j}$$

$$\frac{\partial \hat{y}_{RK}^{re \leftarrow t \rightarrow}}{\partial w_{x_{rs}^{re \leftarrow t \rightarrow}}} = \sum_{r=1}^{n_r} \sum_{s=1}^m \frac{\partial \hat{y}_{RK}^{re \leftarrow t \rightarrow}}{\partial p_{rs}^{re \leftarrow t \rightarrow}} \frac{\partial p_{rs}^{re \leftarrow t \rightarrow}}{\partial w_{x_{rs}^{re \leftarrow t \rightarrow}}} = \sum_{r=1}^{n_r} \sum_{s=1}^m g'_2(p_{rs}^{re \leftarrow t \rightarrow}) \delta_{K,r} \delta_{e,s}$$

$$\cdot s_{js}^{re \leftarrow t \rightarrow} \delta_{r,i} = g'_2(p_{rs}^{re \leftarrow t \rightarrow}) \cdot s_{j \leftarrow t \rightarrow} \delta_{K,i}$$

$$\frac{\partial \hat{y}_{RK}^{re \leftarrow t \rightarrow}}{\partial w_{x_{rs}^{re \leftarrow t \rightarrow}}} = \sum_{r=1}^{n_r} \sum_{s=1}^m \frac{\partial \hat{y}_{RK}^{re \leftarrow t \rightarrow}}{\partial s_{rs}^{re \leftarrow t \rightarrow}} \frac{\partial s_{rs}^{re \leftarrow t \rightarrow}}{\partial w_{x_{rs}^{re \leftarrow t \rightarrow}}} = \sum_{r=1}^{n_r} \sum_{s=1}^m g'_2(p_{rs}^{re \leftarrow t \rightarrow}) \cdot W_{s_{rs}^{re \leftarrow t \rightarrow}} \delta_{e,s}$$

$$\cdot g'_1(z_{rs}^{re \leftarrow t \rightarrow}) \cdot x_{js}^{re \leftarrow t \rightarrow} \delta_{r,i} = g'_2(p_{rs}^{re \leftarrow t \rightarrow}) \cdot W_{s_{rs}^{re \leftarrow t \rightarrow}} g'(z_{je \leftarrow t \rightarrow}) x_{je \leftarrow t \rightarrow}$$

$$\frac{\partial \hat{y}_{RK}^{re \leftarrow t \rightarrow}}{\partial s_{ij}^{re \leftarrow t \rightarrow}} = \sum_{r=1}^{n_r} \sum_{s=1}^m \frac{\partial \hat{y}_{RK}^{re \leftarrow t \rightarrow}}{\partial s_{rs}^{re \leftarrow t \rightarrow}} \cdot \frac{\partial s_{rs}^{re \leftarrow t \rightarrow}}{\partial s_{ij}^{re \leftarrow t \rightarrow}} = \sum_{r=1}^{n_r} \sum_{s=1}^m g'_2(p_{rs}^{re \leftarrow t \rightarrow}) \cdot W_{s_{rs}^{re \leftarrow t \rightarrow}} \delta_{e,s}$$

$$\cdot g'_1(z_{rs}^{re \leftarrow t \rightarrow}) \cdot W_{s_{rs}^{re \leftarrow t \rightarrow}} \cdot \delta_{s,j} =$$

$$= \sum_{r=1}^{n_r} g'_2(p_{rs}^{re \leftarrow t \rightarrow}) \cdot W_{s_{rs}^{re \leftarrow t \rightarrow}} \cdot g'_1(z_{ri}^{re \leftarrow t \rightarrow}) W_{s_{ri}^{re \leftarrow t \rightarrow}} \delta_{e,j}$$

$$\frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{ij} \langle t-1 \rangle} = \sum_{k=1}^{n_s} \sum_{l=1}^m \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{kl} \langle t \rangle} \cdot \frac{\partial S_{kl} \langle t \rangle}{\partial S_{ij} \langle t-1 \rangle} = \sum_{k=1}^{n_s} \sum_{l=1}^m \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{kl} \langle t \rangle} \cdot g'_1(z_{kl} \langle t \rangle) \cdot w_{ski} \delta_{lj}$$

$$= \sum_{k=1}^{n_s} \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{kj} \langle t \rangle} \cdot g'_1(z_{kj} \langle t \rangle) \cdot w_{ski} \rightarrow$$

vectorized implementation  $\rightarrow \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{ij} \langle t-1 \rangle} = w_{st} \cdot \left( \frac{\partial E \langle t+1 : T_y \rangle}{\partial S \langle t \rangle} * g'_1(z \langle t \rangle) \right)$

$$\frac{\partial E \langle t+1 : T_y \rangle}{\partial X_{ij} \langle t \rangle} = w_{sx}^T \left( \frac{\partial E \langle t+1 : T_y \rangle}{\partial S \langle t \rangle} * g'_1(z \langle t \rangle) \right)$$

$$\frac{\partial E \langle t+1 : T_y \rangle}{\partial w_{ski} \langle ij \rangle} = \sum_{k=1}^{n_s} \sum_{l=1}^m \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{kl} \langle t \rangle} \cdot \frac{\partial S_{kl} \langle t \rangle}{\partial w_{ski} \langle ij \rangle} = \sum_{k=1}^{n_s} \sum_{l=1}^m \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{kl} \langle t \rangle} \cdot g'_1(z_{kl} \langle t \rangle) \cdot \delta_{il} \langle t-1 \rangle \cdot \delta_{ki}$$

$$= \sum_{l=1}^m \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{il} \langle t \rangle} \cdot g'_1(z_{il} \langle t \rangle) \cdot \delta_{il} \langle t-1 \rangle$$

vectorized implementation  $\rightarrow \frac{\partial E \langle t+1 : T_y \rangle}{\partial w_{st}} = \left( \frac{\partial E \langle t+1 : T_y \rangle}{\partial S \langle t \rangle} * g'_1(z \langle t \rangle) \right) \cdot S^T \langle t-1 \rangle$

$$\frac{\partial E \langle t+1 : T_y \rangle}{\partial w_{sx}} = \left( \frac{\partial E \langle t+1 : T_y \rangle}{\partial S \langle t \rangle} * g'_1(z \langle t \rangle) \right) \cdot X^T \langle t \rangle$$

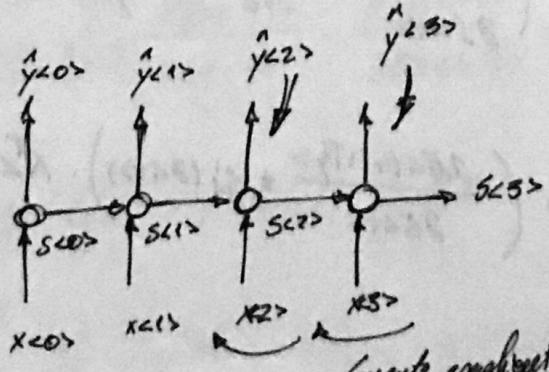
$$\begin{aligned}
 \frac{\partial E^{L+}}{\partial s_{ij}^{L+}} &= \sum_{k=1}^n \sum_{c=1}^m \frac{\partial E^{L+}}{\partial y_{kc}^{L+}} \cdot \frac{\partial y_{kc}^{L+}}{\partial s_{ij}^{L+}} = \sum_{k=1}^n \sum_{c=1}^m \frac{\partial E^{L+}}{\partial y_{kc}^{L+}} \cdot g_2'(p_{kc}^{L+}) \cdot w_{skc} \cdot \delta_{cj} \\
 &= \sum_{k=1}^n \frac{\partial E^{L+}}{\partial y_{kj}^{L+}} \cdot g_2'(p_{kj}^{L+}) \cdot w_{skj} - \\
 \frac{\partial E^{L+}}{\partial w_{skj}} &= \sum_{k=1}^n \sum_{c=1}^m \frac{\partial E^{L+}}{\partial y_{kc}^{L+}} \cdot \frac{\partial y_{kc}^{L+}}{\partial w_{skj}} = \sum_{k=1}^n \left( \sum_{c=1}^m \frac{\partial E^{L+}}{\partial y_{kc}^{L+}} \cdot g_2'(p_{kc}^{L+}) \cdot \delta_{jc} \right) \delta_{ki}
 \end{aligned}$$

Vektorisiert  
Implementation  $\rightarrow \frac{\partial E^{L+}}{\partial s^{L+}} = w_{ps}^T \left( \frac{\partial E^{L+}}{\partial \hat{y}^{L+}} * g_2'(p^{L+}) \right)$

$$\frac{\partial E^{L+}}{\partial w_{ps}} = \left( \frac{\partial E^{L+}}{\partial \hat{y}^{L+}} * g_2'(p^{L+}) \right) \cdot s^T$$

$$\frac{\partial E}{\partial s^{L+1}} = \left[ w_{ps}^T \left( \frac{\partial E^{L+}}{\partial \hat{y}^{L+}} * g_2'(p^{L+}) \right) * j_1'(z^{L+}) \right] + w_{ss}^T \left( \frac{\partial E^{L+1:Tx}}{\partial s^{L+}} * g_1'(z^{L+}) \right)$$

$w_{ss}^T$



Implementation.

compute gradients, start across first & accumulate grad.