

Recursive Deep Models for Semantic Compositionality Over A Sentiment Treebank Derivations for Gradient Equations

Alain Soltani

In the following sections, for a vector b , $b_{[i]}$ denotes the i -th item; for a matrix A , $A_{[i,j]}$, $A_{[i,:]}$, $A_{[:,j]}$ respectively denote the (i, j) coefficient, i -th row and the j -th column.

1 Slices of Tensor V

The RNTN error function for the network's top node p_2 is :

$$E^{p_2} = \sum_{j=1} t_{[j]}^{p_2} \log y_{[j]}^{p_2} = \log y_{[r]}^{p_2}$$

where r is the index of the label (using a 0-1 encoding for the target distribution), so that $t_{[j]}^{p_2} = \mathbb{1}_{\{j=r\}}$.

Thus we can derive the following d -dimensional gradient :

$$\begin{aligned} \frac{\partial E^{p_2}}{\partial p_2} &= \frac{1}{y_{[r]}^{p_2}} \frac{\partial y_{[r]}^{p_2}}{\partial p_2} \\ &= \frac{1}{y_{[r]}^{p_2}} \frac{\partial}{\partial p_2} \left[\frac{\exp(W_{[r,:]}^s p_2 + b_{[r]}^s)}{\sum_{l=1}^d \exp(W_{[l,:]}^s p_2 + b_{[l]}^s)} \right] \\ &= \frac{1}{y_{[r]}^{p_2}} [W_{[r,:]}^s y_{[r]}^{p_2} - y_{[r]}^{p_2} \sum_{l=1}^d W_{[l,:]}^s y_{[l]}^{p_2}], \text{ via a classical product derivation} \\ &= W_{[r,:]}^s - \sum_{l=1}^d W_{[l,:]}^s y_{[l]}^{p_2} \\ &= \sum_{l=1}^d W_{[l,:]}^s t_{[l]}^{p_2} - W_{[l,:]}^s y_{[l]}^{p_2}, \text{ as } t_{[l]}^{p_2} = \mathbb{1}_{\{l=r\}} \\ &= \begin{bmatrix} \sum_{l=1}^d W_{[l,1]}^s (t_{[l]}^{p_2} - y_{[l]}^{p_2}) \\ \vdots \\ \sum_{l=1}^d W_{[l,d]}^s (t_{[l]}^{p_2} - y_{[l]}^{p_2}) \end{bmatrix} = (W^s)^T (t^{p_2} - y^{p_2}). \end{aligned} \tag{1}$$

In the meantime, one can form the following gradient for $V^{[k]}, k \in \langle 1, d \rangle$, k -th slice of tensor V :

$$\frac{\partial p_{2[j]}}{\partial V^{[k]}} = \frac{\partial}{\partial V^{[k]}} f\left(\begin{bmatrix} a \\ p_1 \end{bmatrix}^T V^{[j]} \begin{bmatrix} a \\ p_1 \end{bmatrix} + W_{[j,:]} \begin{bmatrix} a \\ p_1 \end{bmatrix} + b_{[j]}\right), \quad (2)$$

where W, b are respectively the weight and bias matrices.

One can immediately note that $\frac{\partial p_{2[j]}}{\partial V^{[k]}} = 0_{(2d, 2d)}$ if $j \neq k, j \in \langle 1, d \rangle$.

For $j = k$:

$$\begin{aligned} \frac{\partial p_{2[k]}}{\partial V^{[k]}} &= f'\left(\begin{bmatrix} a \\ p_1 \end{bmatrix}^T V^{[k]} \begin{bmatrix} a \\ p_1 \end{bmatrix} + W_{[k,:]} \begin{bmatrix} a \\ p_1 \end{bmatrix} + b_{[k]}\right) \begin{bmatrix} a \\ p_1 \end{bmatrix} \begin{bmatrix} a \\ p_1 \end{bmatrix}^T \\ &= f'(x_{[k]}^{p_2}) \begin{bmatrix} a \\ p_1 \end{bmatrix} \begin{bmatrix} a \\ p_1 \end{bmatrix}^T. \end{aligned} \quad (3)$$

by chain derivation on the function's first term.

Finally, we obtain the following gradient :

$$\begin{aligned} \frac{\partial E^{p_2}}{\partial V^{[k]}} &= \frac{\partial E^{p_2}}{\partial p_2} \frac{\partial p_2}{\partial V^{[k]}} = \frac{\partial E^{p_2}}{\partial p_{2[k]}} \frac{\partial p_{2[k]}}{\partial V^{[k]}} \text{ as there is only one non-zero slice in tensor } \frac{\partial p_2}{\partial V^{[k]}} \\ &= [(W^s)^T (t^{p_2} - y^{p_2}) \otimes f'(x^{p_2})]_{[k]} \begin{bmatrix} a \\ p_1 \end{bmatrix} \begin{bmatrix} a \\ p_1 \end{bmatrix}^T \\ &= \delta_{[k]}^{p_2, com} \begin{bmatrix} a \\ p_1 \end{bmatrix} \begin{bmatrix} a \\ p_1 \end{bmatrix}^T = \delta_{[k]}^{p_2, com} \begin{bmatrix} a \\ p_1 \end{bmatrix} \odot \begin{bmatrix} a \\ p_1 \end{bmatrix}, \end{aligned} \quad (4)$$

\odot denoting the outer product, \otimes the Hadamard product.

2 Weight W

Very similarly as in (2), we can calculate the following gradient, related to W :

$$\begin{aligned} \frac{\partial p_2}{\partial W} &= \frac{\partial}{\partial W} f\left(\begin{bmatrix} a \\ p_1 \end{bmatrix}^T V^{[1:d]} \begin{bmatrix} a \\ p_1 \end{bmatrix} + W \begin{bmatrix} a \\ p_1 \end{bmatrix} + b\right) \\ &= f'\left(\begin{bmatrix} a \\ p_1 \end{bmatrix}^T V^{[1:d]} \begin{bmatrix} a \\ p_1 \end{bmatrix} + W \begin{bmatrix} a \\ p_1 \end{bmatrix} + b\right) \begin{bmatrix} a \\ p_1 \end{bmatrix}^T \\ &= f'(x^{p_2}) \begin{bmatrix} a \\ p_1 \end{bmatrix}^T. \end{aligned} \quad (5)$$

Hence

$$\begin{aligned} \frac{\partial E^{p_2}}{\partial W} &= \frac{\partial E^{p_2}}{\partial p_2} \frac{\partial p_2}{\partial W} = (W^s)^T (t^{p_2} - y^{p_2}) \otimes f'(x^{p_2}) \begin{bmatrix} a \\ p_1 \end{bmatrix}^T \\ &= \delta^{p_2, com} \odot \begin{bmatrix} a \\ p_1 \end{bmatrix}. \end{aligned} \quad (6)$$

3 Bias b

This one is rather straightforward :

$$\begin{aligned}\frac{\partial p_2}{\partial b} &= \frac{\partial}{\partial b} f\left(\begin{bmatrix} a \\ p_1 \end{bmatrix}^T V^{[1:d]} \begin{bmatrix} a \\ p_1 \end{bmatrix} + W \begin{bmatrix} a \\ p_1 \end{bmatrix} + b\right) \\ &= f'\left(\begin{bmatrix} a \\ p_1 \end{bmatrix}^T V^{[1:d]} \begin{bmatrix} a \\ p_1 \end{bmatrix} + W \begin{bmatrix} a \\ p_1 \end{bmatrix} + b\right) = f'(x^{p_2}).\end{aligned}\tag{7}$$

Thus, we obtain by chain derivation

$$\frac{\partial E^{p_2}}{\partial b} = \frac{\partial E^{p_2}}{\partial p_2} \frac{\partial p_2}{\partial b} = (W^s)^T (t^{p_2} - y^{p_2}) \bigotimes f'(x^{p_2}) = \delta^{p_2, com}.\tag{8}$$

4 Classification weight W^s

Let us consider Eq. (1) from Socher's paper for computing the posterior probability over labels :

$$y^a = \text{softmax}(W^s a + b^s).\tag{9}$$

when including the bias b^s . At the network's top node, this becomes :

$$y^{p_2} = \text{softmax}(W^s p_2 + b^s).$$

From this follows the gradient's calculus for coefficient (i, j) from W^s :

$$\frac{\partial y_{[r]}^{p_2}}{\partial W_{[i,j]}^s} = \frac{\partial}{\partial W_{[i,j]}^s} \left[\frac{\exp(W_{[r,:]}^s p_2 + b_{[r]}^s)}{\sum_{l=1}^d \exp(W_{[l,:]}^s p_2 + b_{[l]}^s)} \right]$$

• If $r = i$:

$$\begin{aligned}\frac{\partial y_{[r]}^{p_2}}{\partial W_{[i,j]}^s} &= - \frac{y_{[r]}^{p_2}}{\sum_{l=1}^d \exp(W_{[l,:]}^s p_2 + b_{[l]}^s)} \frac{\partial}{\partial W_{[i,j]}^s} [\exp(W_{[i,:]}^s p_2 + b_{[i]}^s)] \\ &= -y_{[r]}^{p_2} y_{[i]}^{p_2} p_{[j]}^2.\end{aligned}$$

• If $r \neq i$:

$$\frac{\partial y_{[r]}^{p_2}}{\partial W_{[i,j]}^s} = y_{[r]}^{p_2} p_{[j]}^2 - y_{[r]}^{p_2} y_{[i]}^{p_2} p_{[j]}^2.$$

Finally :

$$\frac{\partial E^{p_2}}{\partial W_{[i,j]}^s} = \frac{1}{y_{[r]}^{p_2}} \frac{\partial y_{[r]}^{p_2}}{\partial W_{[i,j]}^s} = p_{[j]}^2 (\mathbb{1}_{\{r=i\}} - y_{[i]}^{p_2}) = p_{[j]}^2 (t_{[i]}^{p_2} - y_{[i]}^{p_2})$$

i.e.

$$\begin{aligned}\frac{\partial E^{p_2}}{\partial W_{[i,j]}^s} &= [(t^{p_2} - y^{p_2}) \bigodot p^2]_{[i,j]} \\ \text{and } \frac{\partial E^{p_2}}{\partial W^s} &= (t^{p_2} - y^{p_2}) \bigodot p^2.\end{aligned}\tag{10}$$

5 Classification bias b^s

Similarly, one can obtain the classification bias' gradient :

$$\begin{aligned} \frac{\partial E^{p_2}}{\partial b_{[i]}^s} &= \frac{1}{y_{[r]}^{p_2}} \frac{\partial y_{[r]}^{p_2}}{\partial b_{[i]}^s} = \frac{1}{y_{[r]}^{p_2}} \frac{\partial}{\partial b_{[i]}^s} \left[\frac{\exp(W_{[r,:]}^s p_2 + b_{[r]}^s)}{\sum_{l=1}^d \exp(W_{[l,:]}^s p_2 + b_{[l]}^s)} \right] \\ &= \mathbb{1}_{\{r=i\}} - y^{p_2}[i] = t^{p_2}[i] - y^{p_2}[i] \end{aligned} \quad (11)$$

and

$$\frac{\partial E^{p_2}}{\partial b^s} = t^{p_2} - y^{p_2}. \quad (12)$$

6 Moving through the nodes

One can then make the calculus for the children's node by first back-propagating the error :

$$\delta^{p_2, down} = (W^T \delta^{p_2, com} + S) \otimes f' \left(\begin{bmatrix} a \\ p_1 \end{bmatrix} \right) \quad (13)$$

This equation is a classical RNN back-propagation error, except for the additional term S accounting for all the tensor slices' contributions :

$$S = \sum_{k=1}^d \delta^{p_2, com} (V^{[k]} + (V^{[k]})^T) \begin{bmatrix} a \\ p_1 \end{bmatrix}. \quad (14)$$

We decompose the stacked error $\delta^{p_2, down}$ in two sub-errors, $\delta_{[1:d]}^{p_2, down}$ and $\delta_{[d+1:2d]}^{p_2, down}$. These are used to compute the children errors :

$$\begin{aligned} \bullet \delta^{p_1, com} &= \delta^{p_1, s} + \delta_{[d+1:2d]}^{p_2, down} \\ \bullet \delta^{a, com} &= \delta_{[1:d]}^{p_2, down} \end{aligned} \quad (15)$$

Thus the full derivative for each parameter (tensor slice, weight, bias, etc.) sums up all node gradients :

$$\begin{aligned} \frac{\partial E}{\partial V^{[k]}} &= \frac{\partial E^{p_2}}{\partial V^{[k]}} + \delta_{[k]}^{p_1, com} \begin{bmatrix} b \\ c \end{bmatrix} \odot \begin{bmatrix} b \\ c \end{bmatrix} \\ \frac{\partial E}{\partial W} &= \frac{\partial E^{p_2}}{\partial W} + \delta^{p_1, com} \odot \begin{bmatrix} b \\ c \end{bmatrix} \\ \frac{\partial E}{\partial b} &= \frac{\partial E^{p_2}}{\partial b} + \delta^{p_1, com} \\ \text{etc.} & \end{aligned} \quad (16)$$