

Machine Learning for Graphs and Sequential Data Exercise Sheet 08

Neural Network Approaches for Sequential Data

Problem 1: Word2vec defines a mapping from a single word to a single fixed vector. Explain and provide an example why this will not be expressive enough regarding homographs (i.e., words with the same spelling but having more than one meaning). Propose an alternative solution.

Word2vec assigns the same vector representation to all words with the same spelling. Hence, it will fail to disambiguate the words that have the same spelling but different meaning, e.g. *duck* – an animal and *duck* – to lower head.

One possible solution is to use RNNs, e.g. LSTMs, to get contextually aware embeddings that will be different for different contexts. Given a sentence, each word will get a hidden state as an embedding. The training task is to predict the next word based on all the previous. Another example is attention based models, e.g., BERT, GPT-3 etc.

Problem 2: Given a previous hidden state $\mathbf{h}^{(t-1)} \in \mathbb{R}^D$ and a current input $\mathbf{x}^{(t)} \in \mathbb{R}^N$, the recurrent neural network equations to update the hidden state and produce the output are:

$$\mathbf{z}^{(t)} = \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}$$

$$\mathbf{h}^{(t)} = \tanh\left(\mathbf{z}^{(t)}\right)$$

$$\mathbf{o}^{(t)} = \mathbf{V}\mathbf{h}^{(t)}$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}\left(\mathbf{o}^{(t)}\right)$$

where parameters $\mathbf{W} \in \mathbb{R}^{D \times D}$, $\mathbf{U} \in \mathbb{R}^{D \times N}$ and $\mathbf{V} \in \mathbb{R}^{M \times D}$ are shared at every step.

To train an RNN we need gradients of loss w.r.t. the parameters: $\partial L / \partial \mathbf{W}$, $\partial L / \partial \mathbf{U}$ and $\partial L / \partial \mathbf{V}$. Your task is to arrive at the equations given on slide 17 in the lecture.

Use the fact that $\partial L / \partial \mathbf{o}^{(t)} = \hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}$, where $\mathbf{y}^{(t)}$ is the true output.

Hint: Since parameters are shared, the total gradient is the sum of the contributions over all the steps. Because of that, it might be easier to introduce copies of parameters, e.g. $\mathbf{W}^{(t)}$ – a copy of \mathbf{W} at step t , calculate $\partial L / \partial \mathbf{W}^{(t)}$ and sum over all t .

We want to show the following formula:

$$\frac{\partial L}{\partial \mathbf{V}} = \sum_t \left(\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)} \right) \left(\mathbf{h}^{(t)} \right)^T \quad (1)$$

$$\frac{\partial L}{\partial \mathbf{W}} = \sum_t \text{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^2 \right) \frac{\partial L}{\partial \mathbf{h}^{(t)}} \left(\mathbf{h}^{(t-1)} \right)^T \quad (2)$$

$$\frac{\partial L}{\partial \mathbf{U}} = \sum_t \text{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^2 \right) \frac{\partial L}{\partial \mathbf{h}^{(t)}} \left(\mathbf{x}^{(t)} \right)^T \quad (3)$$

We start with (1). We can unroll the gradient into all the contributions across time, and then apply the chain rule:

$$\frac{\partial L}{\partial \mathbf{V}} = \sum_t \frac{\partial L^{(t)}}{\partial \mathbf{V}} = \sum_t \frac{\partial L^{(t)}}{\partial \mathbf{o}^{(t)}} \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{V}}$$

Each element of this product can be written as:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{o}^{(t)}} &= \hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}, \\ \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{V}} &= \left(\mathbf{h}^{(t)} \right)^T, \end{aligned}$$

where we slightly abuse the notation in the second term. The value $\frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{V}}$ is actually a tensor, but most of the entries are zero. See also: <https://web.stanford.edu/class/cs224n/readings/gradient-notes.pdf>.

We now consider (2). Again, expand the different time contributions and apply the chain rule:

$$\frac{\partial L}{\partial \mathbf{W}} = \sum_t \frac{\partial L}{\partial \mathbf{W}^{(t)}} = \sum_t \frac{\partial L}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{W}^{(t)}}$$

The Jacobian $\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{z}^{(t)}}$ is a diagonal matrix since $h_i^{(t)}$ does not depend on $z_j^{(t)}$ if $i \neq j$. The tanh operation is applied element-wise. Thus, $\frac{\partial h_i^{(t)}}{\partial z_i^{(t)}} = 1 - \tanh(z_i^{(t)})^2 = 1 - (\mathbf{h}_i^{(t)})^2$. Hence, we have:

$$\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{z}^{(t)}} = \text{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^2 \right),$$

and also:

$$\frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{W}^{(t)}} = \left(\mathbf{h}^{(t-1)} \right)^T.$$

For (3), computations are similar except that we use $\frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{U}^{(t)}} (\mathbf{x}^{(t-1)})^T$ instead of $\frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{W}^{(t)}} = (\mathbf{h}^{(t-1)})^T$. These shows the formula at slide 17.

Additionally, note that we could compute the derivative $\frac{\partial L}{\partial \mathbf{h}^{(t)}}$ recursively from $\frac{\partial L}{\partial \mathbf{h}^{(t+1)}}$:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{h}^{(t)}} &= \sum_{s=t}^T \frac{\partial L^{(s)}}{\partial \mathbf{h}^{(t)}} \\ &= \frac{\partial L^{(t)}}{\partial \mathbf{h}^{(t)}} + \sum_{s=t+1}^T \frac{\partial L^{(s)}}{\partial \mathbf{h}^{(t)}} \end{aligned}$$

We can compute the derivatives:

$$\frac{\partial L^{(t)}}{\partial \mathbf{h}^{(t)}} = \frac{\partial L^{(t)}}{\partial \mathbf{o}^{(t)}} \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}} = \mathbf{V}^T (\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)})$$
$$\sum_{s=t+1}^T \frac{\partial L^{(s)}}{\partial \mathbf{h}^{(t)}} = \sum_{s=t+1}^T \frac{\partial L^{(s)}}{\partial \mathbf{h}^{(t+1)}} \frac{\partial \mathbf{h}^{(t+1)}}{\partial \mathbf{z}^{(t+1)}} \frac{\partial \mathbf{z}^{(t+1)}}{\partial \mathbf{h}^{(t)}} = \sum_{s=t+1}^T \mathbf{W}^T \text{diag} \left(1 - \left(\mathbf{h}^{(t+1)} \right)^2 \right) \frac{\partial L^{(s)}}{\partial \mathbf{h}^{(t+1)}}$$

As explained in slide 18, we see that the gradient depends on future times.

Problem 3: What do you need to change in the equations that you got in the previous exercise if the output $\mathbf{o}^{(t)}$ is used as an input to another neural network?

Instead of using $(\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)})$ we would need to calculate $\partial L / \partial \mathbf{o}^{(t)}$ depending on the subsequent layers and loss. Everything else remains the same.
