

Language modeling is a central task in NLP, and language models can be found at the heart of speech recognition, machine translation, and many other systems. Given a sequence of words (represented as one-hot row vectors)  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$ , a language model predicts the next word  $\mathbf{x}^{(t+1)}$  by modeling:

$$P(\mathbf{x}^{(t+1)} = \mathbf{v}_j \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})$$

where  $\mathbf{v}_j$  is a word in the vocabulary.

Your job is to compute the gradients of a recurrent neural network language model, which uses feedback information in the hidden layer to model the “history”  $\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(1)}$ . Formally, the model<sup>6</sup> is, for  $t = 1, \dots, n - 1$ :

$$\begin{aligned} \mathbf{e}^{(t)} &= \mathbf{x}^{(t)} \mathbf{L} \\ \mathbf{h}^{(t)} &= \text{sigmoid} \left( \mathbf{h}^{(t-1)} \mathbf{H} + \mathbf{e}^{(t)} \mathbf{I} + \mathbf{b}_1 \right) \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax} \left( \mathbf{h}^{(t)} \mathbf{U} + \mathbf{b}_2 \right) \\ \bar{P}(\mathbf{x}^{(t+1)} = \mathbf{v}_j \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) &= \hat{y}_j^{(t)} \end{aligned}$$

where  $\mathbf{h}^{(0)} = \mathbf{h}_0 \in \mathbb{R}^{D_h}$  is some initialization vector for the hidden layer and  $\mathbf{x}^{(t)} \mathbf{L}$  is the product of  $\mathbf{L}$  with the one-hot row vector  $\mathbf{x}^{(t)}$  representing the current word. The parameters are:

$$\mathbf{L} \in \mathbb{R}^{|V| \times d} \quad \mathbf{H} \in \mathbb{R}^{D_h \times D_h} \quad \mathbf{I} \in \mathbb{R}^{d \times D_h} \quad \mathbf{b}_1 \in \mathbb{R}^{D_h} \quad \mathbf{U} \in \mathbb{R}^{D_h \times |V|} \quad \mathbf{b}_2 \in \mathbb{R}^{|V|} \quad (1)$$

where  $\mathbf{L}$  is the embedding matrix,  $\mathbf{I}$  the input word representation matrix,  $\mathbf{H}$  the hidden transformation matrix, and  $\mathbf{U}$  is the output word representation matrix.  $\mathbf{b}_1$  and  $\mathbf{b}_2$  are biases.  $d$  is the embedding dimension,  $|V|$  is the vocabulary size, and  $D_h$  is the hidden layer dimension.

The output vector  $\hat{\mathbf{y}}^{(t)} \in \mathbb{R}^{|V|}$  is a probability distribution over the vocabulary. The model is trained by minimizing the (un-regularized) cross-entropy loss:

$$J^{(t)}(\theta) = CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \sum_{j=1}^{|V|} y_j^{(t)} \log \hat{y}_j^{(t)}$$

where  $\mathbf{y}^{(t)}$  is the one-hot vector corresponding to the target word (which here is equal to  $\mathbf{x}^{(t+1)}$ ). We average the cross-entropy loss across all examples (i.e., words) in a sequence to get the loss for a single sequence.

(a) **(written)** Conventionally, when reporting performance of a language model, we evaluate on *perplexity*, which is defined as:

$$\text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = \frac{1}{\bar{P}(\mathbf{x}_{\text{pred}}^{(t+1)} = \mathbf{x}^{(t+1)} \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})} = \frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} \cdot \hat{y}_j^{(t)}}$$

i.e. the inverse probability of the correct word, according to the model distribution  $\bar{P}$ . Show how you can derive perplexity from the cross-entropy loss (*Hint: remember that  $\mathbf{y}^{(t)}$  is one-hot!*), and thus argue that minimizing the (arithmetic) mean cross-entropy loss will also minimize the (geometric) mean perplexity across the training set. ***This should be a very short problem - not too perplexing!***

For a vocabulary of  $|V|$  words, what would you expect perplexity to be if your model predictions were completely random (chosen uniformly from the vocabulary)? Compute the corresponding cross-entropy loss for  $|V| = 10000$ .

- (b) **(written)** Compute the gradients of the loss  $J$  with respect to the following model parameters at a single point in time  $t$  (to save a bit of time, you don't have to compute the gradients with the respect to  $\mathbf{U}$  and  $\mathbf{b}_1$ ):

$$\frac{\partial J^{(t)}}{\partial \mathbf{b}_2} \quad \frac{\partial J^{(t)}}{\partial \mathbf{L}_{x^{(t)}}} \quad \frac{\partial J^{(t)}}{\partial \mathbf{I}} \Big|_{(t)} \quad \frac{\partial J^{(t)}}{\partial \mathbf{H}} \Big|_{(t)}$$

where  $\mathbf{L}_{x^{(t)}}$  is the row of  $\mathbf{L}$  corresponding to the current word  $\mathbf{x}^{(t)}$ , and  $\Big|_{(t)}$  denotes the gradient for the appearance of that parameter at time  $t$  (equivalently,  $\mathbf{h}^{(t-1)}$  is taken to be fixed, and you need not backpropagate to earlier timesteps just yet - you'll do that in part (c)).

Additionally, compute the derivative with respect to the *previous* hidden layer value:

$$\frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}}$$

- (c) **(coding)** Use the corpus provided in the file to train a language model, using a sequence of words to predict the following word. The basic structure is already finished, complete the model by filling the code in ***model.py*** and ***train&test.py***.