

The architecture is a RNN language model. The model is formally defined as:

$$\begin{aligned} e^{(t)} &= x^{(t)} L \\ h^{(t)} &= \text{sigmoid} \left(h^{(t-1)} H + e^{(t)} I + b_1 \right) \\ \hat{y}^{(t)} &= \text{softmax} \left(h^{(t)} U + b_2 \right) \end{aligned}$$

$$p(x_{t+1} = v_j | x_t, \dots, x_1) = \hat{y}_j^{(t)}$$

where L is the embedding matrix, I the input word representation matrix, H the hidden transformation matrix, and U is the output word representation matrix. b1 and b2 are biases. d is the embedding dimension, |V| is the vocabulary size, and Dh is the hidden layer dimension.

2. Compute the gradients for all the model parameters at a single point in time t:

$$\frac{\partial J^{(t)}}{\partial U} \quad \frac{\partial J^{(t)}}{\partial b_2} \quad \frac{\partial J^{(t)}}{\partial L_{x^{(t)}}} \quad \frac{\partial J^{(t)}}{\partial I} \Big|_{(t)} \quad \frac{\partial J^{(t)}}{\partial H} \Big|_{(t)} \quad \frac{\partial J^{(t)}}{\partial b_1} \Big|_{(t)}$$

The solution is provided as:

The partial derivatives:

$$\begin{aligned} \frac{\partial J^{(t)}}{\partial U} &= \left(\mathbf{h}^{(t)} \right)^T (\mathbf{y} - \hat{\mathbf{y}}) \\ \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t)}} &= (\mathbf{y} - \hat{\mathbf{y}}) \mathbf{U}^T \\ \frac{\partial J^{(t)}}{\partial b_2} &= (\mathbf{y} - \hat{\mathbf{y}}) \\ \frac{\partial J^{(t)}}{\partial b_1} \Big|_{(t)} &= \left(\frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t)}} \odot \text{sigmoid}' \left(\mathbf{h}^{(t-1)} \mathbf{H} + e^{(t)} \mathbf{I} + b_1 \right) \right) \\ \frac{\partial J^{(t)}}{\partial L_{x^{(t)}}} &= \frac{\partial J^{(t)}}{\partial e^{(t)}} = \frac{\partial J^{(t)}}{\partial b_1} \Big|_{(t)} I^T \\ \frac{\partial J^{(t)}}{\partial \mathbf{I}} \Big|_{(t)} &= \left(e^{(t)} \right)^T \frac{\partial J^{(t)}}{\partial b_1} \Big|_{(t)} \\ \frac{\partial J^{(t)}}{\partial \mathbf{H}} \Big|_{(t)} &= \left(h^{(t-1)} \right)^T \frac{\partial J^{(t)}}{\partial b_1} \Big|_{(t)} \\ \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}} &= \mathbf{H}^T \frac{\partial J^{(t)}}{\partial b_1} \Big|_{(t)} \end{aligned}$$

TODO: I was wondering (taking the first partial derivative as an example) where the h-transpose comes from in the calculation for the first partial derivative.