

METL

TP2

Tientso Ning

1. Perplexity and corresponding CE loss is defined as:

$$PP^{(t)}(y^{(t)}, \hat{y}^{(t)}) = \frac{1}{p(x_{t+1}^{pred}=x_{t+1}|x_t, \dots, x_1)} = \frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} \cdot \hat{y}_j^{(t)}}$$
$$= -\sum_{i=1}^{|V|} y_i^{(t)} \cdot \log \hat{y}_i^{(t)}$$

Since we have a one hot representation, there will be only one non-zero slot.

$$= -y_i(t) \cdot \log(\hat{y}_i(t)) = -\log(\hat{y}_i(t))$$
$$= \log(1) - \log\left(\hat{y}_i^{(t)}\right) = \log\left(\frac{1}{\hat{y}_i(t)}\right), \text{ when } Pp = \frac{1}{\hat{y}(t)}$$

Therefore $CE(y^{(t)}, \hat{y}^{(t)}) = \log PP^{(t)}(y^{(t)}, \hat{y}^{(t)})$

For $|V| = 2000$, and $|V| = 10,000$ $\log(2000) = 7.600902459542082$ $\log(10000) = 9.210340371976184$

For a vocab of $|V|$ words, complete randomness would mean all have the same probabilities $1/|V|$. Therefore Perplexity will be $|V|$.

2. Derive the gradients.

$$\frac{\partial J^{(t)}}{\partial U} = (h^{(t)})^\top (\hat{y} - y)$$

Where $h^{(t)}$ is the hidden-layer output at time t, and $\hat{y} - y$ is our cost function derivative.

$$\frac{\partial J^{(t)}}{\partial h^{(t)}} = (\hat{y} - y)U^T \text{ Since } U \text{ is the output layer matrix.}$$

$\frac{\partial J^{(t)}}{\partial b_2} = (\hat{y} - y)$ The value for the bias “flattens out” and we don’t have any other values since the bias is added at the last layer.

$$\left. \frac{\partial J^{(t)}}{\partial b_i} \right|_t = \left. \frac{\partial J^{(t)}}{\partial h^{(t)}} \odot \theta' \left(h^{(t-1)}H + e^{(t)}I + b \right) \right|_t$$

The values in the sigmoid are $h^{(t-1)}H$, which is the hidden layer, $e^{(t)}I$ corresponds to the embedding matrix of the input representation, with bias b1.

$$\frac{\partial J^{(4)}}{\partial L_x^{(H)}} = \frac{\partial J^{(t)}}{\partial e^{(t)}} = \left. \frac{\partial J^{(t)}}{\partial b_1} \right|_t \cdot I^\top$$

from the values inside the sigmoid above.

$$\left. \frac{\partial J^{(t)}}{\partial I} \right|_t = (e^{(t)})^\top \cdot \left. \frac{\partial J^{(t)}}{\partial b_1} \right|_t$$

same as above, but for the other term in the sigmoid.

$$\left. \frac{\partial J^{(t)}}{\partial tt} \right|_t = (h^{(t-1)})^\top \left. \frac{\partial J^{(t)}}{\partial b_i} \right|_t \frac{\partial J^{(t)}}{\partial h^{(t-1)}} = H^T \left. \frac{\partial J^{(t)}}{\partial b_1} \right|_t$$

3. Draw an unrolled RNN and derive gradients.

$\delta^{(t-1)} = \frac{\partial J^{(t)}}{\partial h^{(t-1)}}$ is the error term.

$$\left. \frac{\partial J^{(t)}}{\partial b_1} \right|_{t-1} = \frac{\partial J^{(t)}}{\partial h^{(t-1)}} \odot \theta' \left(h^{(t-2)}H + e^{(t-1)}I + b_1 \right)$$
$$= \delta^{(t-1)} \odot \theta' \left(h^{(t-2)}H + e^{(t-1)}I + b_1 \right)$$

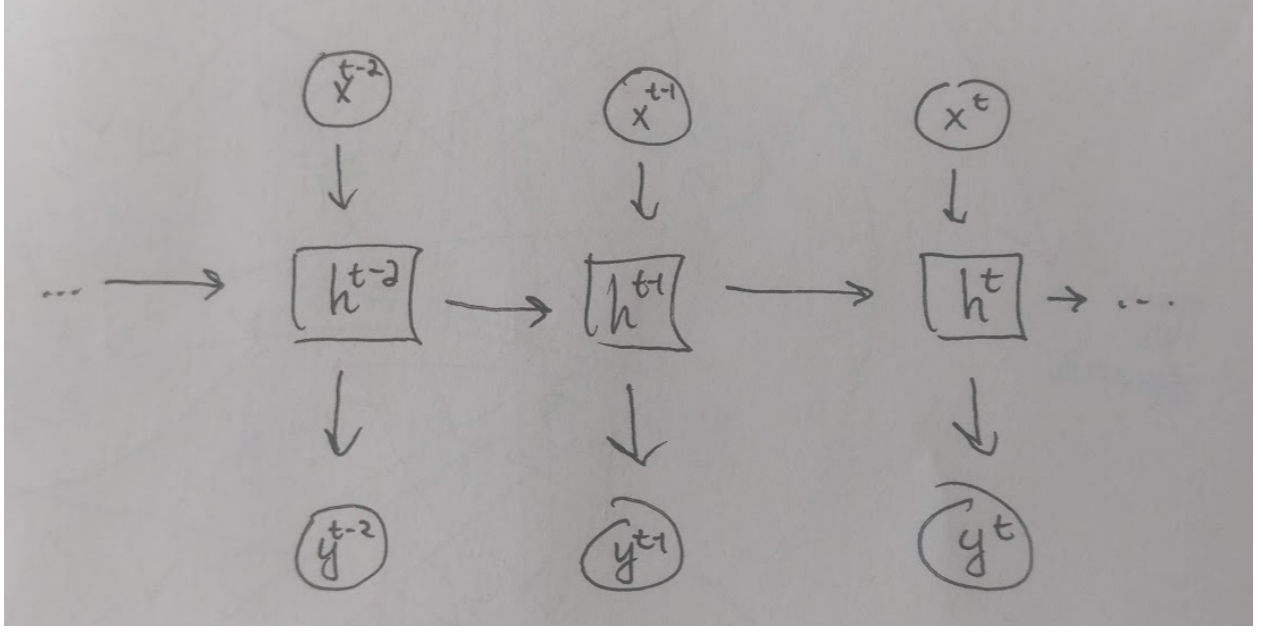


Figure 1: unrolled

from the sigmoid application, which is an element-wise application, at the layers. Where the values $h^{t-2}H$ from the hidden layer, $e^{(t-1)}I$ is from the embedding representation with b1 bias.

$$\frac{\partial J^{(t)}}{\partial L_x^{(t-1)}} = \frac{\partial J^{(t)}}{\partial b_i} \Big|_{t-1} I^T$$

Since $L_x^{(t-1)}$ is the column corresponding to the word x at time $t-1$, and $e^t = x^t L$ where I is the input word representation.

$$\frac{\partial J^{(t)}}{\partial I} \Big|_{t-1} = (e^{(t-1)})^T \cdot \frac{\partial J^{(t)}}{\partial b_1} \Big|_{t-1}$$

$$\frac{\partial J^{(t)}}{\partial (H)} \Big|_{t-1} = (h^{t-2})^T \cdot \frac{\partial J^{(t)}}{\partial b_1} \Big|_{t-1}$$

both from the terms in the sigmoid above.

4. Define the complexities.

$$h^{(t-1)} = d(D_h) \quad e^{(t)} = O(d) \quad \hat{y}^{(t)} = O(|V|)$$

e^t is $x_t L$, meaning $O(d)$ since d is the embedding dimension. L is size d . $h^{(t)}$ and \hat{y}^t are derived from derivatives. Forwards and backwards have the same time complexity and z is just a scalar to complexity.

The slow step is calculating \hat{y}^t since the application of the softmax would require the dimension to be of $|V|$, and the vocabulary sizes can be big. We can be working in the magnitudes of hundreds or thousands of words.