

CS 224n : Assignment 2

3. Recurrent Neural Networks Language Modelling

a) Perplexity

$$p p^{(t)}(y^{(t)}, \hat{y}^{(t)})$$

$$= \frac{1}{\bar{p}(x_{pred}^{(t+1)} | x^{(t)} \dots x^{(1)})}$$

$$= \frac{1}{\sum_{j=1}^V y_j^{(t)} \cdot \hat{y}_j^{(t)}}$$

We know that $y^{(t)}$ is one-hot:

$$y^{(t)} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow k = (\delta_{ik})$$

$$p p^{(t)}(y^{(t)}, \hat{y}^{(t)}) = \frac{1}{1 \cdot \hat{y}_k^{(t)}} \\ = \frac{1}{\hat{y}_k^{(t)}}$$

Cross-entropy loss

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

$$\begin{aligned} &= 0 \quad \forall j \neq k \\ &= 1 \quad \text{if } j = k \end{aligned}$$

$$= - \log \hat{y}_k^{(t)}$$

$$= - \log \frac{1}{p p^{(t)}}$$

$$J^{(t)}(\theta) = \log p p^{(t)}$$

$$pp^{(n)}(y^{(n)}, \hat{y}^{(n)}) = \exp(J^{(n)}(\theta))$$

a)ii) We want to:

$$\text{minimize} \left[\frac{1}{T} \sum_{t=1}^T pp(y^{(t)}, \hat{y}^{(t)}) \right]^{1/T}$$

$$\text{minimize} \frac{1}{T} \log \left(\frac{1}{T} \sum_{t=1}^T pp(y^{(t)}, \hat{y}^{(t)}) \right) \quad \downarrow \log$$

$$\text{minimize} \frac{1}{T} \sum_{t=1}^T \log pp(y^{(t)}, \hat{y}^{(t)}) \quad \downarrow =$$

$$\text{minimize} \frac{1}{T} \sum_{t=1}^T CE(y^{(t)}, \hat{y}^{(t)}) \quad \downarrow =$$

a)iii)

Let's pick the next word uniformly.

$$pp^{(n)}(y^{(n)}, \hat{y}^{(n)}) = \frac{1}{P(x_{pred}^{(n+1)} \neq x^{(n+1)} | x^{(1)} \dots x^{(n)})}$$

$$= \frac{1}{\frac{1}{|V|}} = |V|$$

Then, the corresponding cross-entropy loss would be:

$$\begin{aligned} J^{(n)}(\theta) &= \log pp^{(n)} \\ &= \log |V| \\ &= \log(10^4) \end{aligned}$$

$$J^{(n)}(\theta) = 4$$

b).

let's define the following intermediary terms:

$$z^{(t)} = W_h h^{(t-1)} + W_e e^{(t)} + b_1 \quad (1)$$

$$\text{so that: } h^{(t)} = \sigma(z^{(t)}) \quad (2)$$

$$g^{(t)} = U \cdot h^{(t)} + b_2 \quad (3)$$

$$\hat{y}^{(t)} = \text{softmax}(g^{(t)}) \quad (4)$$

$$J^{(t)}(\theta) = \text{CE}(y^{(t)}, \hat{y}^{(t)})$$

$$\boxed{\frac{\partial J^{(t)}}{\partial U}} \xrightarrow{\substack{\in \mathbb{R} \\ \in \mathbb{R}^{1 \times 1} \times \mathbb{R}^{1 \times D_h}}} \frac{\partial J^{(t)}}{\partial U} = \frac{\partial J^{(t)}}{\partial \theta^{(t)}} \frac{\partial \theta^{(t)}}{\partial U} = \underbrace{\hat{y}^{(t)} - y^{(t)}}_{\substack{\text{see assignment 1)}}} = h^{(t)} \quad (D_h, 1)$$

$$\boxed{\frac{\partial J^{(t)}}{\partial U} = \delta_1^{(t)} \cdot h^{(t)T} \quad \text{where } \delta_1^{(t)} = \hat{y}^{(t)} - y^{(t)}} \quad \in (1 \times 1, 1)$$

$$\frac{\partial J^{(t)}}{\partial \theta^{(t)}} = \frac{\partial J^{(t)}}{\partial \theta^{(t)}} \frac{\partial \theta^{(t)}}{\partial h^{(t)}} \frac{\partial h^{(t)}}{\partial z^{(t)}} \frac{\partial z^{(t)}}{\partial e^{(t)}} \quad (d, 1)$$

$$= (\hat{y} - y)^T \cdot U \cdot \underbrace{\sigma'(z^{(t)})}_{\text{element-wise multiplication, as } \sigma \text{ is a scalar function}} \cdot W_e \quad (1 \times 1, 1) \quad (1 \times 1, D_h) \quad (D_h, D_h) \quad (D_h, d)$$

$$= \underbrace{(\hat{y} - y)^T U \circ \sigma'(z^{(t)})}_{(1, d)} W_e$$

should be equal to its transpose

$$\boxed{\frac{\partial J^{(t)}}{\partial e^{(t)}} = \underbrace{(\sigma'(z^{(t)}) W_e)^T}_{(d, 1)}}$$

$$\begin{aligned} \text{where } \delta_2^{(t)} &= (\hat{y} - y)^T U \circ \sigma'(z^{(t)}) \\ &= \sigma(z^{(t)}) \circ (1 - \sigma(z^{(t)})) \\ &= h^{(t)} \circ (1 - h^{(t)}) \end{aligned}$$

$$\boxed{\delta_2^{(t)} = \underbrace{(\hat{y} - y)^T}_{(1, 1 \times 1)} \underbrace{U \circ h^{(t)} \circ (1 - h^{(t)})}_{(1 \times 1, D_h) \quad (1 \times 1, D_h) \quad (D_h, D_h)}}$$

$$\left. \frac{\partial J(t)}{\partial w_e} \right|_t = \frac{\partial J}{\partial z(t)} \cdot \frac{\partial z(t)}{\partial w_e}$$

$$\underbrace{(D_n, 1)}_{(1, D_n)} = \underbrace{\delta_2(t)}_{(1, D_n)} \cdot \underbrace{e(t)}_{(D_n, 1)}$$

$$\left. \frac{\partial J(t)}{\partial w_e} \right|_t = \delta_2(t)^T e(t)^T$$

$$\left. \frac{\partial J(t)}{\partial w_h} \right|_t = \frac{\partial J}{\partial z(t)} \cdot \frac{\partial z(t)}{\partial w_h}$$

$$\underbrace{(D_n, D_n)}_{(1, D_n)} = \underbrace{\delta_2(t)}_{(1, D_n)} \cdot \underbrace{h^{(t-1)}}_{(D_n, 1)}$$

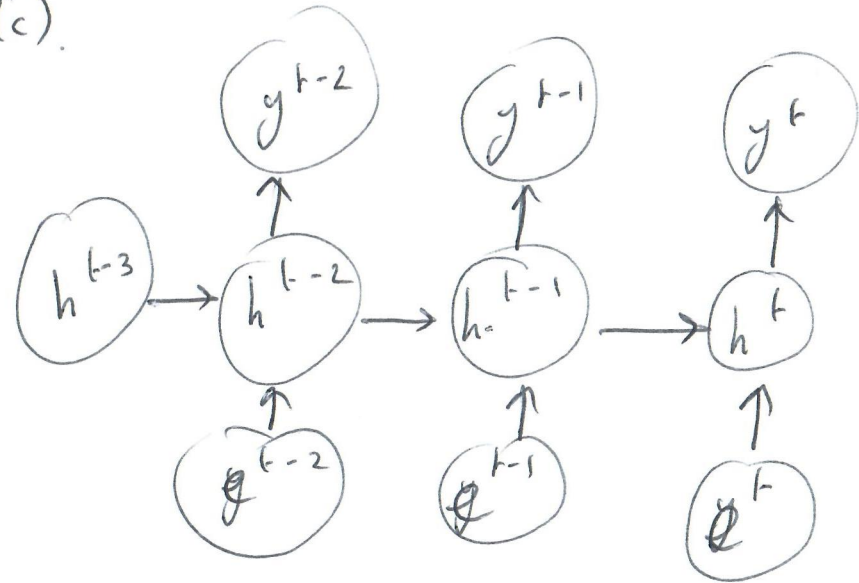
$$\left. \frac{\partial J(t)}{\partial w_h} \right|_t = \delta_2(t)^T h^{(t-1)T}$$

$$\frac{\partial J}{\partial h^{(t-1)}} = \frac{\partial J}{\partial z(t)} \cdot \frac{\partial z(t)}{\partial h^{(t-1)}}$$

$$\underbrace{(D_n, 1)}_{(1, D_h)} = \underbrace{\delta_2(t)}_{(1, D_h)} \cdot \underbrace{w_h}_{(D_h, D_n)}$$

$$\frac{\partial J}{\partial h^{(t-1)}} = (\delta_2(t) w_h)^T$$

(c).



$$\frac{\partial J(t)}{\partial e^{(t-1)}} = \frac{\partial J(t)}{\partial h^{(t-1)}} \cdot \frac{\partial h^{(t-1)}}{\partial z^{(t-1)}} \cdot \frac{\partial z^{(t-1)}}{\partial e^{(t-1)}}$$

$d, 1$ $\gamma^{(t-1)}$ $\circ h^{(t-1)} \circ (1-h^{(t-1)})$ W_e

$D_n, 1$ D_n, D_n D_n, d

$$\frac{\partial J(t)}{\partial e^{(t-1)}} = \left(\gamma^{(t-1)} W_e \right)^T \circ h^{(t-1)} \circ (1-h^{(t-1)})$$

$$\frac{\partial J(t)}{\partial W_e} \Big|_{t-1} = \frac{\partial J(t)}{\partial h^{(t-1)}} \cdot \frac{\partial h^{(t-1)}}{\partial z^{(t-1)}} \cdot \frac{\partial z^{(t-1)}}{\partial W_e} \Big|_{t-1}$$

(D_n, d) $\gamma^{(t-1)}$ $\circ h^{(t-1)} \circ (1-h^{(t-1)})$ $e^{(t-1)}$

$(D_n, 1)$ (D_n, D_n) $(d, 1)$

$$\frac{\partial J(t)}{\partial W_e} \Big|_{t-1} = \gamma^{(t-1)T} \circ h^{(t-1)} \circ (1-h^{(t-1)}) e^{(t-1)T}$$

$$\frac{\partial J(t)}{\partial W_h} \Big|_{t-1} = \frac{\partial J(t)}{\partial h^{(t-1)}} \cdot \frac{\partial h^{(t-1)}}{\partial z^{(t-1)}} \cdot \frac{\partial z^{(t-1)}}{\partial W_h} \Big|_{t-1}$$

(D_n, D_n) $\gamma^{(t-1)}$ $\circ h^{(t-1)} \circ (1-h^{(t-1)})$ $h^{(t-2)}$

$(D_n, 1)$ (D_n, D_n) $(D_n, 1)$

$$\frac{\partial J(t)}{\partial W_h} \Big|_{t-1} = \gamma^{(t-1)T} \circ h^{(t-1)} \circ (1-h^{(t-1)}) h^{(t-2)T}$$

(d) Calculating complexity

$\frac{\partial J(t)}{\partial v} : D_n \cdot |V|$ operations

$\frac{\partial J(t)}{\partial e^{(t)}} : |V| \cdot D_n + D_n \cdot d$ operations

$\frac{\partial J(t)}{\partial W_e} \Big|_t : D_n \cdot d$ operations

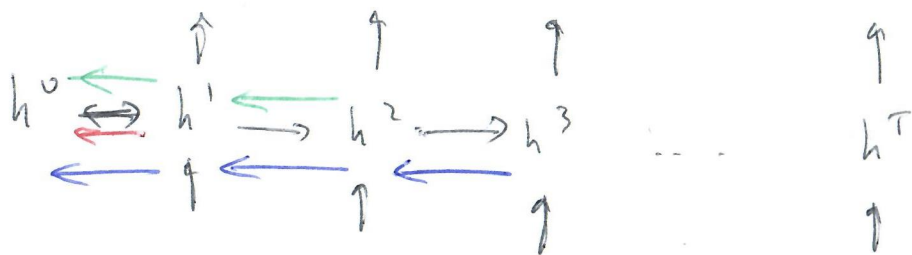
$\frac{\partial J(t)}{\partial W_h} \Big|_t : D_n \cdot D_n$ operations

$\frac{\partial J(t)}{\partial h^{(t-1)}} : D_n \cdot D_n$ operations

Which makes a total of:

$$\mathcal{O}(D_h \cdot |V| + D_h \cdot d + D_h^2) \text{ operations}$$

e) We could backpropagate all our losses at the same time, which would make us do:



1 pass for h^1

2 for h^2

\vdots
 τ for h^T

$$\Sigma = \frac{\tau(\tau+1)}{2} = \mathcal{O}(\tau^2)$$

But we could also wait for those passes coming from the ~~left~~ right to all be computed, before "packaging" them together and sending them on the left.

→ aggregating flux into 1.

Hence, we'd need just τ passes.

Complexity for Twosds

$$= \mathcal{O}(\tau D_h (|V| + d + D_h))$$

f). $|V|$ = vocabulary : $\sim 10^4 - 6$

d = dimension of our embeddings : $\sim 10^2$

D_h : dimension of hidden layer : $\sim 10^2$

$$\Rightarrow \mathcal{O}(\tau \cdot D_h \cdot |V|)$$