

CS 224n: Deep learning for NLP

Assignment 2.

2) 3)

Stack	buffer	new dependency	transition
[root]	[I, parsed, this, sentence, correctly]		Initial Config.
[root, I]	[Parsed, this, sentence, correctly]		SHIFT
[root, I, Parsed]	[this, sentence, correctly]		SHIFT
[root, Parsed]	[this, sentence, correctly]	Parsed \rightarrow I	Left-Arc
[root, Parsed, this]	[sentence, correctly]		SHIFT
[root, Parsed, this, sentence]	[correctly]		SHIFT
[root, Parsed, sentence]	[correctly]	Sentence \rightarrow This	Left-Arc
[root, Parsed]	[correctly]	Parsed \rightarrow sentence	Right-Arc
[root, Parsed, Correctly]	[]		Shift
[root, Parsed]	[]	Parsed \rightarrow Correctly	Right-Arc
[root]	[]	Root \rightarrow Parsed	Right-Arc

g(i) ~~keeping~~ the momentum accumulated throughout the descent of the gradient can help us continue update our parameters in case

- the "real" gradient becomes null
- When we get stuck at a local minimum for instance
- It could help against vanishing gradient

g) (ii) dividing by \sqrt{t} has the effect of increasing the final gradient of previously low values:
ie: ~~very~~ very low values of gradient now have slightly bigger values:
 \Rightarrow this may help combat vanishing gradient

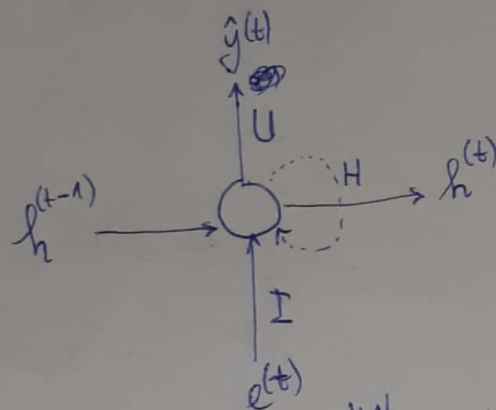
$$3) P(x^{(t+1)} = y_j | x^{(t)} \dots x^{(1)}) = \hat{y}_j^{(t)}$$

$$e^{(t)} = x^{(t)} \cdot L$$

$$h^{(t)} = \text{sigmoid}(h^{(t-1)} H + e^{(t)} \cdot I + b_1) = \text{sigmoid}(\theta^{(t)})$$

$$\hat{y}^{(t)} = \text{softmax}(h^{(t)} U + b_2) = \text{softmax}(z^{(t)})$$

$$h^{(0)} = h_0 \in \mathbb{R}^{D_h}, L \in \mathbb{R}^{|V| \times d}, H \in \mathbb{R}^{D_h \times D_h}, I \in \mathbb{R}^{d \times D_h}, U \in \mathbb{R}^{D_h \times |V|}, b_2 \in \mathbb{R}^{|V|}$$



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

$$J_{\text{seq}}(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta)$$

$$(2) PP^{(t)}(y^{(t)}, \hat{y}^{(t)}) = \frac{1}{P(x_{\text{pred}}^{(t+1)} = x^{(t+1)} | x^{(1)} \dots x^{(t)})} = \frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} \cdot \hat{y}_j^{(t)}}$$

(Perplexity)

$$J_{\text{seq}} = -\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^{|V|} y_j^{(t)} \log(\hat{y}_j^{(t)})$$

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

$$= \log\left(\frac{1}{\hat{y}_k^{(t)}}\right) = \log\left(\frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} \hat{y}_j^{(t)}}\right) = \log(PP^{(t)})$$

$$J^{(t)} = \log(PP^{(t)}) \iff PP^{(t)} = 2^{J^{(t)}}$$

(2) * let show that:

minimizing $J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}$ \Leftrightarrow minimizing $PP(\theta) = \left[\prod_{t=1}^T (PP^{(t)}) \right]^{\frac{1}{T}}$

$T = |V|$

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)} = \frac{1}{T} \sum_{t=1}^T \log(PP^{(t)}(\theta))$$

$$= \frac{1}{T} \log \left[\prod_{t=1}^T PP^{(t)}(\theta) \right]$$

$$= \log \left[\left(\prod_{t=1}^T PP^{(t)}(\theta) \right)^{\frac{1}{T}} \right]$$

$$\Rightarrow \underset{\theta}{\operatorname{argmin}} J(\theta) = \underset{\theta}{\operatorname{argmin}} \log \left[\left(\prod_{t=1}^T PP^{(t)}(\theta) \right)^{\frac{1}{T}} \right]$$

because log is strictly \nearrow

$$= \underset{\theta}{\operatorname{argmin}} \left[\left(\prod_{t=1}^T PP^{(t)}(\theta) \right)^{\frac{1}{T}} \right]$$

* $PP^{(t)}(\theta) = \frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} \hat{y}_j^{(t)}}$; if predictions are random from uniform distribution

$\rightarrow \forall j \hat{y}_j^{(t)} = \frac{1}{|V|}$

$$\Rightarrow PP^{(t)} = \frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} \cdot \frac{1}{|V|}} = \frac{1}{y_k^{(t)} \cdot \frac{1}{|V|}} \quad \left| \text{Where } k \text{ is the correct class} \right.$$

$$PP^{(t)} = \frac{1}{\frac{1}{|V|}} = |V|$$

in the case $|V| = 10000$

$$J^{(t)} = \log(PP^{(t)}) = \log(10000)$$

$$J^{(t)} \approx 9.210$$

3) b) let $z^{(t)} = h^{(t)} u + b_2$
 $\theta^{(t)} = h^{(t-1)} H + e^{(t)} I + b_1$

$$\frac{\partial J^{(t)}}{\partial z^{(t)}} = g^{(t)} - y^{(t)} = \delta_1^{(t)}$$

$$\frac{\partial J^{(t)}}{\partial \theta^{(t)}} = \frac{\partial J^{(t)}}{\partial z^{(t)}} \cdot \frac{\partial z^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial \theta^{(t)}} = \delta_1^{(t)} \cdot u^T \odot \sigma(\theta^{(t)}) \odot (1 - \sigma(\theta^{(t)})) = \delta_2^{(t)}$$

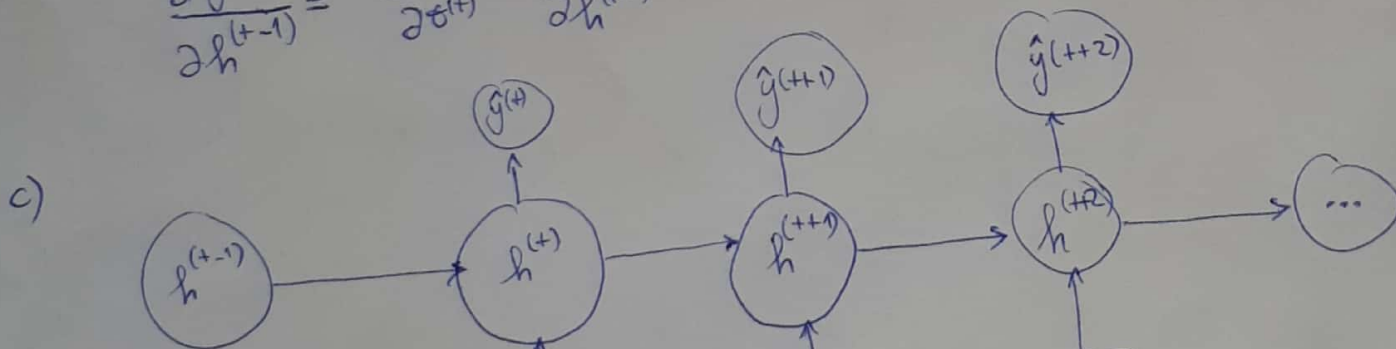
$$\frac{\partial J^{(t)}}{\partial b_2} = \frac{\partial J^{(t)}}{\partial z^{(t)}} \cdot \frac{\partial z^{(t)}}{\partial b_2} = \delta_1^{(t)} \cdot \frac{\partial}{\partial b_2} (h^{(t)} \cdot u + b_2) = \delta_1^{(t)}$$

$$\frac{\partial J^{(t)}}{\partial L_{x^{(t)}}} = \frac{\partial J^{(t)}}{\partial e^{(t)}} = \frac{\partial J^{(t)}}{\partial \theta^{(t)}} \cdot \frac{\partial \theta^{(t)}}{\partial e^{(t)}} = \delta_2^{(t)} \cdot I^T$$

$$\left. \frac{\partial J^{(t)}}{\partial I} \right|_{(t)} = \frac{\partial J^{(t)}}{\partial \theta^{(t)}} \cdot \left. \frac{\partial \theta^{(t)}}{\partial I} \right|_{(t)} = \delta_2^{(t)} \cdot (e^{(t)})^T$$

$$\left. \frac{\partial J^{(t)}}{\partial H} \right|_{(t)} = \frac{\partial J^{(t)}}{\partial \theta^{(t)}} \cdot \left. \frac{\partial \theta^{(t)}}{\partial H} \right|_{(t)} = \delta_2^{(t)} \cdot (h^{(t-1)})^T$$

$$\frac{\partial J^{(t)}}{\partial h^{(t-1)}} = \frac{\partial J^{(t)}}{\partial \theta^{(t)}} \cdot \frac{\partial \theta^{(t)}}{\partial h^{(t-1)}} = \delta_2^{(t)} \cdot H^T = \delta^{(t-1)}$$



$$\frac{\partial J^{(t)}}{\partial L_{x^{(t-1)}}} = \frac{\partial J^{(t)}}{\partial e^{(t-1)}} = \frac{\partial J^{(t)}}{\partial \theta^{(t)}} \cdot \frac{\partial \theta^{(t)}}{\partial e^{(t-1)}} = \delta_2^{(t)} \cdot \frac{\partial \theta^{(t)}}{\partial h^{(t-1)}} \cdot \frac{\partial h^{(t-1)}}{\partial e^{(t-1)}} = \delta^{(t-1)} \cdot I^T \cdot \sigma'(\theta^{(t-1)})$$

$$\left. \frac{\partial J^{(t)}}{\partial I} \right|_{t-1} = \sum_{k=t-1}^t \frac{\partial J^{(k)}}{\partial h^{(k)}} \cdot \frac{\partial h^{(k)}}{\partial I} = \frac{\partial J^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial I} + \frac{\partial J^{(t)}}{\partial h^{(t-1)}} \cdot \frac{\partial h^{(t-1)}}{\partial I}$$

$$\Rightarrow \left[\frac{\partial J^{(t)}}{\partial I} \right]_{t-1} = \delta_2^{(t)} (e^{(t)})^T + \delta^{(t-1)} \odot \sigma'(\theta^{(t-1)}) e^{(t-1)T}$$

$$\begin{aligned}
 \frac{\partial \mathcal{F}^{(t)}}{\partial \mathbf{H}} \Big|_{(t-1)} &= \sum_{k=t-1}^t \frac{\partial \mathcal{F}^{(k)}}{\partial \mathbf{h}^{(k)}} \cdot \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{H}} \\
 &= \frac{\partial \mathcal{F}^{(t)}}{\partial \mathbf{h}^{(t)}} \cdot \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{H}} + \frac{\partial \mathcal{F}^{(t-1)}}{\partial \mathbf{h}^{(t-1)}} \cdot \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{H}} \\
 &= \frac{\partial \mathcal{F}^{(t)}}{\partial \mathbf{H}} + \delta^{(t-1)} \odot \sigma'(\mathbf{z}^{(t-1)}) \cdot (\mathbf{h}^{(t-2)})^T
 \end{aligned}$$

$$\boxed{\frac{\partial \mathcal{F}^{(t)}}{\partial \mathbf{H}} \Big|_{(t-1)} = \delta_z^{(t)} (\mathbf{h}^{(t-1)})^T + \delta^{(t-1)} \odot \sigma'(\mathbf{z}^{(t-1)}) \cdot (\mathbf{h}^{(t-2)})^T}$$

d)

$$\begin{aligned}
 \# \text{operations}(\mathcal{F}^{(t)} | \mathbf{h}^{(t-1)}) &= O(\text{op}(\hat{\mathbf{y}}^{(t)})) = O(\text{op}(\mathbf{h}^{(t)}) \times M \times D_h^2 + |V|) \\
 \# \text{op}(\mathbf{h}^{(t)}) &= O(\#(\mathbf{h}^{(t-1)}) \times D_h^3) + \dots
 \end{aligned}$$