

2. (a)

stack	buffer	new dependency	transition
[ROOT]	[I, parsed, this, sentence, correctly]		Initial Configuration
[ROOT, I]	[parsed, this, sentence, correctly]		SHIFT
[ROOT, I, parsed]	[this, sentence, correctly]		SHIFT
[ROOT, parsed]	[this, sentence, correctly]	parsed → I	LEFT-ARC
[ROOT, parsed, this]	[sentence, correctly]		SHIFT
[ROOT, parsed, this, sentence]	[correctly]		SHIFT
[ROOT, parsed, sentence]	[correctly]	sentence → this	LEFT-ARC
[ROOT, parsed]	[correctly]	parsed → sentence	RIGHT-ARC
[ROOT, parsed, correctly]	∅		SHIFT
[ROOT, parsed]	∅	parsed → correctly	RIGHT-ARC
[ROOT]	∅		RIGHT-ARC

2. (b) A sentence containing n words will be parsed in $2n$ steps, because each word needs to be pushed into the stack, which needs n steps, and then popped out of the stack until only ROOT is in the stack, which needs another n steps.

3. (a) (i) To answer this question, we need to take into consideration that only one element of $\{\mathbf{y}\}$ is 1, and others 0. Assume that the i th element is 1. Then for the cross-entropy, we have

$$J = -y_i \ln \hat{y}_i = -\ln \hat{y}_i.$$

Correspondingly, the expression for the perplexity becomes

$$PP = \frac{1}{\hat{y}_i}.$$

Therefore,

$$e^J = e^{-\ln \hat{y}_i} = e^{\ln \hat{y}_i^{-1}} = \frac{1}{\hat{y}_i} = PP.$$

3. (a) (ii) Apply logarithm to the geometric mean perplexity:

$$\begin{aligned} \ln \left(\prod_{t=1}^T PP^{(t)} \right)^{1/T} &= \frac{1}{T} \left(\ln PP^{(1)} + \ln PP^{(2)} + \dots + \ln PP^{(T)} \right) \\ &= \frac{1}{T} \left(J^{(1)} + J^{(2)} + \dots + J^{(T)} \right) = \frac{1}{T} \sum_{t=1}^T J^{(t)}. \end{aligned}$$

Since the logarithm function is an increasing one, minimizing the geometric mean perplexity is equivalent to minimizing the arithmetic mean cross-entropy loss.

3. (a) (iii) At any step, the probability of the model predicting the correct word is $\frac{1}{|V|}$. So the perplexity is $|V| = 100000$. The cross-entropy loss is $\ln |V| = 4 \ln 10 \approx 9.21$.

3. (b) We already know that $\frac{\partial J}{\partial \{\boldsymbol{\theta}^{(t)}\}} = \{\hat{\mathbf{y}}^{(t)}\} - \{\mathbf{y}^{(t)}\}$, therefore

$$\begin{aligned}
\frac{\partial J^{(t)}}{\partial [\mathbf{U}]} &= \left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right) \frac{\partial \{\boldsymbol{\theta}\}}{\partial [\mathbf{U}]} \\
&= \begin{bmatrix} \left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right) \begin{bmatrix} \frac{\partial \theta_1}{\partial \langle \mathbf{u}_1 \rangle} \\ \vdots \\ \frac{\partial \theta_1}{\partial \langle \mathbf{u}_V \rangle} \\ \frac{\partial \theta_2}{\partial \langle \mathbf{u}_1 \rangle} \\ \vdots \\ \frac{\partial \theta_2}{\partial \langle \mathbf{u}_V \rangle} \end{bmatrix} \\ \vdots \\ \left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right) \begin{bmatrix} \frac{\partial \theta_V}{\partial \langle \mathbf{u}_1 \rangle} \\ \vdots \\ \frac{\partial \theta_V}{\partial \langle \mathbf{u}_V \rangle} \end{bmatrix} \end{bmatrix} \\
&= \begin{bmatrix} \left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right) \begin{bmatrix} \langle \mathbf{h}^{(t)} \rangle \\ \langle \mathbf{0} \rangle \\ \vdots \\ \langle \mathbf{0} \rangle \\ \langle \mathbf{0} \rangle \\ \langle \mathbf{h}^{(t)} \rangle \\ \vdots \\ \langle \mathbf{0} \rangle \end{bmatrix} \\ \vdots \\ \left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right) \begin{bmatrix} \langle \mathbf{0} \rangle \\ \vdots \\ \langle \mathbf{h}^{(t)} \rangle \end{bmatrix} \end{bmatrix} \\
&= \begin{bmatrix} \left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right)_{*1} \cdot \langle \mathbf{h}^{(t)} \rangle \\ \left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right)_{*2} \cdot \langle \mathbf{h}^{(t)} \rangle \\ \vdots \\ \left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right)_{*V} \cdot \langle \mathbf{h}^{(t)} \rangle \end{bmatrix} \\
&= \left\{ \mathbf{h}^{(t)} \right\} \left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right).
\end{aligned}$$

$$\frac{\partial \{\boldsymbol{\theta}^{(t)}\}}{\partial \langle \mathbf{h}^{(t)} \rangle} = [\mathbf{U}],$$

$$\frac{\partial \{\mathbf{h}^{(t)}\}}{\partial \langle \mathbf{z}^{(t)} \rangle} = \text{diag} \left(\langle \mathbf{h}^{(t)} \rangle \odot \left(\langle \mathbf{1} \rangle - \langle \mathbf{h}^{(t)} \rangle \right) \right).$$

$$\frac{\partial \{\mathbf{z}^{(t)}\}}{\partial \langle \mathbf{e}^{(t)} \rangle} = [\mathbf{W}_e].$$

Therefore,

$$\begin{aligned} \frac{\partial J^{(t)}}{\partial \langle \mathbf{e}^{(t)} \rangle} &= \frac{\partial J^{(t)}}{\partial \langle \boldsymbol{\theta}^{(t)} \rangle} \frac{\partial \{\boldsymbol{\theta}^{(t)}\}}{\partial \langle \mathbf{h}^{(t)} \rangle} \frac{\partial \{\mathbf{h}^{(t)}\}}{\partial \langle \mathbf{z}^{(t)} \rangle} \frac{\partial \{\mathbf{z}^{(t)}\}}{\partial \langle \mathbf{e}^{(t)} \rangle} \\ &= \left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right) [\mathbf{U}] \text{diag} \left(\langle \mathbf{h}^{(t)} \rangle \odot \left(\langle \mathbf{1} \rangle - \langle \mathbf{h}^{(t)} \rangle \right) \right) [\mathbf{W}_e]. \end{aligned}$$

$$\frac{\partial J^{(t)}}{\partial \{\mathbf{e}^{(t)}\}} = [\mathbf{W}_e]^\top \text{diag} \left(\left(\{\mathbf{1}\} - \{\mathbf{h}^{(t)}\} \right) \odot \{\mathbf{h}^{(t)}\} \right) [\mathbf{U}]^\top \left(\{\hat{\mathbf{y}}^{(t)}\} - \{\mathbf{y}^{(t)}\} \right).$$

$$\begin{aligned} \frac{\partial \{\mathbf{z}^{(t)}\}}{\partial [\mathbf{W}_e]} &= \begin{bmatrix} \left[\begin{array}{c} \frac{\partial z_1}{\partial \langle \mathbf{W}_{e,1*} \rangle} \\ \vdots \\ \frac{\partial z_1}{\partial \langle \mathbf{W}_{e,D_h*} \rangle} \end{array} \right] \\ \left[\begin{array}{c} \frac{\partial z_2}{\partial \langle \mathbf{W}_{e,1*} \rangle} \\ \vdots \\ \frac{\partial z_2}{\partial \langle \mathbf{W}_{e,D_h*} \rangle} \end{array} \right] \\ \vdots \\ \left[\begin{array}{c} \frac{\partial z_{D_h}}{\partial \langle \mathbf{W}_{e,1*} \rangle} \\ \vdots \\ \frac{\partial z_{D_h}}{\partial \langle \mathbf{W}_{e,D_h*} \rangle} \end{array} \right] \end{bmatrix} \\ &= \begin{bmatrix} \left[\begin{array}{c} \langle \mathbf{e}^{(t)} \rangle \\ \langle \mathbf{0} \rangle \\ \vdots \\ \langle \mathbf{0} \rangle \end{array} \right] \\ \left[\begin{array}{c} \langle \mathbf{0} \rangle \\ \langle \mathbf{e}^{(t)} \rangle \\ \vdots \\ \langle \mathbf{0} \rangle \end{array} \right] \\ \vdots \\ \left[\begin{array}{c} \langle \mathbf{0} \rangle \\ \vdots \\ \langle \mathbf{e}^{(t)} \rangle \end{array} \right] \end{bmatrix}. \end{aligned}$$

$$\begin{aligned}\left.\frac{\partial J^{(t)}}{\partial [\mathbf{W}_e]}\right|_t &= \frac{\partial J^{(t)}}{\partial \langle \boldsymbol{\theta}^{(t)} \rangle} \frac{\partial \{\boldsymbol{\theta}^{(t)}\}}{\partial \langle \mathbf{h}^{(t)} \rangle} \frac{\partial \{\mathbf{h}^{(t)}\}}{\partial \langle \mathbf{z}^{(t)} \rangle} \frac{\partial \{\mathbf{z}^{(t)}\}}{\partial [\mathbf{W}_e]} \\ &= \left(\left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right) [\mathbf{U}] \odot \langle \mathbf{h}^{(t)} \rangle \odot \left(\langle \mathbf{1} \rangle - \langle \mathbf{h}^{(t)} \rangle \right) \right)^\top \otimes \langle \mathbf{e}^{(t)} \rangle.\end{aligned}$$

Similarly,

$$\left.\frac{\partial J^{(t)}}{\partial [\mathbf{W}_h]}\right|_t = \left(\left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right) [\mathbf{U}] \odot \langle \mathbf{h}^{(t)} \rangle \odot \left(\langle \mathbf{1} \rangle - \langle \mathbf{h}^{(t)} \rangle \right) \right)^\top \otimes \langle \mathbf{h}^{(t-1)} \rangle.$$

$$\frac{\partial J^{(t)}}{\partial \{\mathbf{h}^{(t-1)}\}} = [\mathbf{W}_h]^\top \text{diag} \left(\left(\langle \mathbf{1} \rangle - \langle \mathbf{h}^{(t)} \rangle \right) \odot \langle \mathbf{h}^{(t)} \rangle \right) [\mathbf{U}]^\top \left(\langle \hat{\mathbf{y}}^{(t)} \rangle - \langle \mathbf{y}^{(t)} \rangle \right).$$

3. c

$$\begin{aligned}\frac{\partial J^{(t)}}{\partial \langle \mathbf{e}^{(t-1)} \rangle} &= \frac{\partial J^{(t)}}{\partial \langle \mathbf{h}^{(t-1)} \rangle} \frac{\partial \{\mathbf{h}^{(t-1)}\}}{\partial \langle \mathbf{z}^{(t-1)} \rangle} \frac{\partial \{\mathbf{z}^{(t-1)}\}}{\partial \langle \mathbf{e}^{(t-1)} \rangle} \\ &= \langle \boldsymbol{\gamma}^{(t-1)} \rangle \text{diag} \left(\langle \mathbf{h}^{(t-1)} \rangle \odot \left(\langle \mathbf{1} \rangle - \langle \mathbf{h}^{(t-1)} \rangle \right) \right) [\mathbf{W}_e].\end{aligned}$$

$$\begin{aligned}\left.\frac{\partial J^{(t)}}{\partial [\mathbf{W}_e]}\right|_{t-1} &= \frac{\partial J^{(t)}}{\partial \langle \mathbf{h}^{(t-1)} \rangle} \frac{\partial \{\mathbf{h}^{(t-1)}\}}{\partial \langle \mathbf{z}^{(t-1)} \rangle} \frac{\partial \{\mathbf{z}^{(t-1)}\}}{\partial [\mathbf{W}_e]} \\ &= \left(\langle \boldsymbol{\gamma}^{(t-1)} \rangle \odot \langle \mathbf{h}^{(t-1)} \rangle \odot \left(\langle \mathbf{1} \rangle - \langle \mathbf{h}^{(t-1)} \rangle \right) \right)^\top \otimes \langle \mathbf{e}^{(t-1)} \rangle.\end{aligned}$$

$$\left.\frac{\partial J^{(t)}}{\partial [\mathbf{W}_h]}\right|_{t-1} = \left(\langle \boldsymbol{\gamma}^{(t-1)} \rangle \odot \langle \mathbf{h}^{(t-1)} \rangle \odot \left(\langle \mathbf{1} \rangle - \langle \mathbf{h}^{(t-1)} \rangle \right) \right)^\top \otimes \langle \mathbf{h}^{(t-2)} \rangle.$$

3. d $\mathcal{O}(|V| D_h + dD_h + D_h^2)$.

3. e $\mathcal{O}(T(|V| D_h + dD_h + D_h^2))$.

3. f $\mathcal{O}(|V| D_h)$ would be the dominant part.