

① Word-level neural Bidirectional Language Model

(a)

$$\left[\nabla_{\theta} \in (\underline{y}, \hat{y}) \right]; \quad i = 2 \quad \frac{d}{d \theta_2} \left[-\log \left[\frac{\exp(\theta_2)}{\sum_j \exp(\theta_j)} \right] \right] =$$

\uparrow

Assume: $(\underline{y})_2 = 1$

$$= - \left[\frac{\sum_j \exp(\theta_j)}{\exp(\theta_2)} \right] \cdot \left[\frac{\exp(\theta_2) \left[\sum_j \exp(\theta_j) - \exp(\theta_2) \right]}{\left[\sum_j \exp(\theta_j) \right]^2} \right] =$$

$$= - \left[\frac{\sum_j \exp(\theta_j) - \exp(\theta_2)}{\sum_j \exp(\theta_j)} \right] = \frac{\exp(\theta_2)}{\sum_j \exp(\theta_j)} - 1 =$$

$= \text{softmax}(\theta_2) - 1$

$$\dots \stackrel{i \neq 2}{=} \frac{d}{d \theta_i} \left[-\log \left[\frac{\exp(\theta_i)}{\sum_j \exp(\theta_j)} \right] \right] =$$

$$= \frac{\sum_j \exp(\theta_j)}{\exp(\theta_i)} \cdot \left[\frac{\exp(\theta_i) \cdot \exp(\theta_j)}{\left(\sum_j \exp(\theta_j) \right)^2} \right] =$$

$$= \frac{\exp(\theta_i)}{\sum_j \exp(\theta_j)} \quad \boxed{= \text{softmax}(\theta_i)}.$$

① ↴

= 7

$$\boxed{\nabla_{\underline{y}} C(\underline{y}, \underline{y}) = \underline{y} - \underline{y}}$$

[a]

(b)

Denote $\underline{\theta} = \underline{h}W_2 + \underline{b}_2 \Rightarrow \hat{\underline{y}} = \text{softmax}(\underline{\theta})$.

$\underline{r} = \underline{x}W_1 + \underline{b}_1 \Rightarrow \underline{h} = \sigma(\underline{r})$

$$\frac{\partial J}{\partial x} = \left(\frac{\partial J}{\partial \underline{\theta}} \right) \cdot \left(\frac{\partial \underline{\theta}}{\partial \underline{h}} \right) \cdot \left(\frac{\partial \underline{h}}{\partial \underline{r}} \right) \cdot \left(\frac{\partial \underline{r}}{\partial x} \right)$$

diagonal-Jacobian:
 $\begin{bmatrix} \sigma(r)_1 \cdot \sigma'(r)_1 \\ \vdots \\ \sigma(r)_h \cdot \sigma'(r)_h \end{bmatrix}$

Chain rule:

$$= [\hat{\underline{y}} - \underline{y}] \left[W_2^T \right] \left[\begin{array}{c|c|c} & M & \end{array} \right] \left[W_1^T \right] =$$

$$[\hat{\underline{y}} - \underline{y}] \left[W_2^T \right] \left[\text{diag}[h \circ (1-h)] \right] \cdot \left[W_1^T \right].$$

[b]

Perplexity: 112.8171402875714

2)

[2] THEORETICAL INQUIRY OF SIMPLE RNN LANGUAGE-MODEL

• denote: $\underline{\theta}^t = h^t V + b_2$, $\underline{y}^t - \underline{\hat{y}}^t = \underline{\mu}^t$; $V^t = H^{t-1}H + e^t I + b_1$

$$(a) \frac{\partial J^{(t)}}{\partial V_{ij}} = \left(\frac{\partial J^t}{\partial \underline{\theta}^t} \right) \left(\frac{\partial \underline{\theta}^t}{\partial V_{ij}} \right) =$$

$$= [\underline{\hat{y}}^t - \underline{\hat{y}}^t] \begin{bmatrix} \vdots \\ h_i^t \\ \vdots \\ 0 \end{bmatrix} \xleftarrow{\text{j'th row}} = [\underline{\hat{y}}^t - \underline{\hat{y}}^t] \cdot h^t;$$

$$\Rightarrow \boxed{\frac{\partial J^{(t)}}{\partial V_i}} = \begin{bmatrix} \underline{\mu}_1^t \cdot h_1^t & \dots & \underline{\mu}_m^t \cdot h_m^t \\ \underline{\mu}_1^t \cdot h_2^t & \dots & \vdots \\ \vdots & \dots & \vdots \\ \underline{\mu}_1^t \cdot h_n^t & \dots & \underline{\mu}_m^t \cdot h_n^t \end{bmatrix} = [h^t]^T \cdot [\underline{\hat{y}}^t - \underline{\hat{y}}^t] =$$

Outer-product.

$$\boxed{[h^t] \cdot [\underline{\hat{y}}^t - \underline{\hat{y}}^t]} = (h^t)^T \cdot [\underline{\hat{y}}^t - \underline{\hat{y}}^t]$$

Outer-product

[when viewed as row vectors!].

$$\left(\frac{\partial J^t}{\partial \underline{\theta}^t} \right) = [\underline{\hat{y}}^t - \underline{\hat{y}}^t]; \quad \frac{\partial J^t}{\partial h^t} = \left(\frac{\partial J^t}{\partial \underline{\theta}^t} \right) \left(\frac{\partial \underline{\theta}^t}{\partial h^t} \right) \left(\frac{\partial h^t}{\partial h^t} \right) =$$

$$= [\underline{\hat{y}}^t - \underline{\hat{y}}^t] V^T \cdot \text{diag}[h^t \circ (1-h^t)]$$

↑
Hadamard-Product.

①

$$\begin{aligned}
 \frac{\partial J^t}{\partial b_{2i}} &= \left(\frac{\partial J^t}{\partial \theta_i} \right) \left(\frac{\partial \theta_i}{\partial b_{2i}} \right) = \\
 &= [\hat{y}^t - \hat{y}^e] \cdot \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \leftarrow i^{\text{th}} \text{ row} = [\underline{\mu}^t]_i
 \end{aligned}$$

$$\Rightarrow \boxed{\frac{\partial J^t}{\partial b_2}} = [\underline{\mu}^t]_1 \dots [\underline{\mu}^t]_m = \underline{\mu}^t = \boxed{[\hat{y}^t - \hat{y}^e]}$$

2 ↴

$$\boxed{\frac{\partial J^t}{\partial b_1}} = \left(\frac{\partial J^t}{\partial v^t} \right) \left(\frac{\partial v^t}{\partial b_1} \right) = \left(\frac{\partial J^t}{\partial v^t} \right) \cdot I_{D \times D}$$

$$= \boxed{[\hat{y}^t - y^t] V^T \cdot \text{diag}[h^t_0 (1-h^t)]}.$$

$$\boxed{\frac{\partial J^t}{\partial I}} = \left(\frac{\partial J^t}{\partial v^t} \right) \left(\frac{\partial v^t}{\partial I} \right) = \boxed{[(C^t)^T \cdot [\hat{y}^t - y^t] V^T \cdot \text{diag}[h^t_0 (1-h^t)]}$$

$$\boxed{\frac{\partial J^t}{\partial H}} = \left(\frac{\partial J^t}{\partial v^t} \right) \left(\frac{\partial v^t}{\partial H} \right) = \boxed{(h^{in})^T \cdot [\hat{y}^t - y^t] V^T \cdot \text{diag}[h^t_0 (1-h^t)]}$$

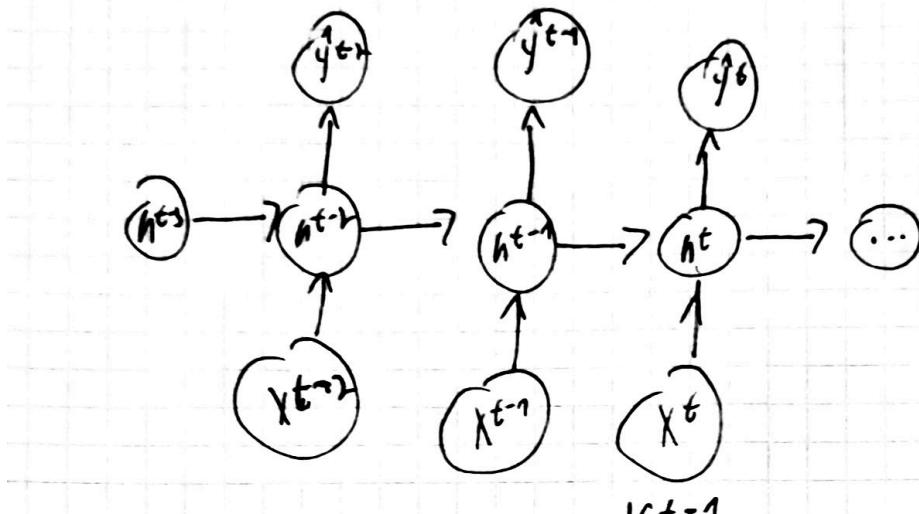
$$\boxed{\frac{\partial J^t}{\partial L_{x^t}}} = \left(\frac{\partial J^t}{\partial v^t} \right) \left(\frac{\partial v^t}{\partial e^t} \right) \left(\frac{\partial e^t}{\partial L_{x^t}} \right) = \boxed{[\hat{y}^t - y^t] V^T \cdot \text{diag}[h^t_0 (1-h^t)] \cdot I^T}$$

$$\boxed{\frac{\partial J^t}{\partial h^{t+1}}} = \left(\frac{\partial J^t}{\partial v^t} \right) \left(\frac{\partial v^t}{\partial h^{t+1}} \right) = \boxed{[\hat{y}^t - y^t] V^T \cdot \text{diag}[h^t_0 (1-h^t)] \cdot H^T} = \boxed{f^{t+1}}$$

(a)

(3)

(b)



$$h^{t-1} = \sigma(h^{t-2}H + C^{t-1}I + b_1)$$

$$h^t = \sigma(h^{t-1}H + C^tI + b_1)$$

$$\left. \left(\frac{\partial J^t}{\partial I} \right) \right|_{t-1} = \left(\frac{\partial J^t}{\partial h^{t-1}} \right) \left[\left(\frac{\partial h^{t-1}}{\partial V^{t-1}} \right) \left(\frac{\partial V^{t-1}}{\partial I} \right) \right] =$$

$$= [C^{t-1}]^T [f^{t-1}] [\text{diag}[h^{t-1} \circ (1 - h^{t-1})]].$$

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

$$= [h^{t-2}]^T [f^{t-1}] [\text{diag}[h^{t-1} \circ (1 - h^{t-1})]]$$



$$\left(\frac{\partial J^t}{\partial Lx^{t-1}} \right) \Big|_{t-1} = \left(\frac{\partial J^t}{\partial h^{t-1}} \right) \cdot \left(\frac{\partial h^{t-1}}{\partial R^{t-1}} \right) \cdot \left(\frac{\partial R^{t-1}}{\partial C^{t-1}} \right) \cdot \left(\frac{\partial C^{t-1}}{\partial Lx^{t-1}} \right) =$$

$$= (f^{t-1}) \cdot (\text{diag}[h^{t-1} \circ (1-h^{t-1})]) \cdot (I)^T.$$

$$\left(\frac{\partial J^t}{\partial b_1} \right) \Big|_{t-1} = \left(\frac{\partial J^t}{\partial h^{t-1}} \right) \left(\frac{\partial h^{t-1}}{\partial r^{t-1}} \right) \left(\frac{\partial r^{t-1}}{\partial b_1} \right) =$$

$$= (f^{t-1}) (\text{diag}[\sigma(r^{t-1}) \circ \sigma(-r^{t-1})]) (I_{bH \times D_H}) =$$

$$= (f^{t-1}) [\text{diag}[h^{t-1} \circ (1-h^{t-1})]].$$

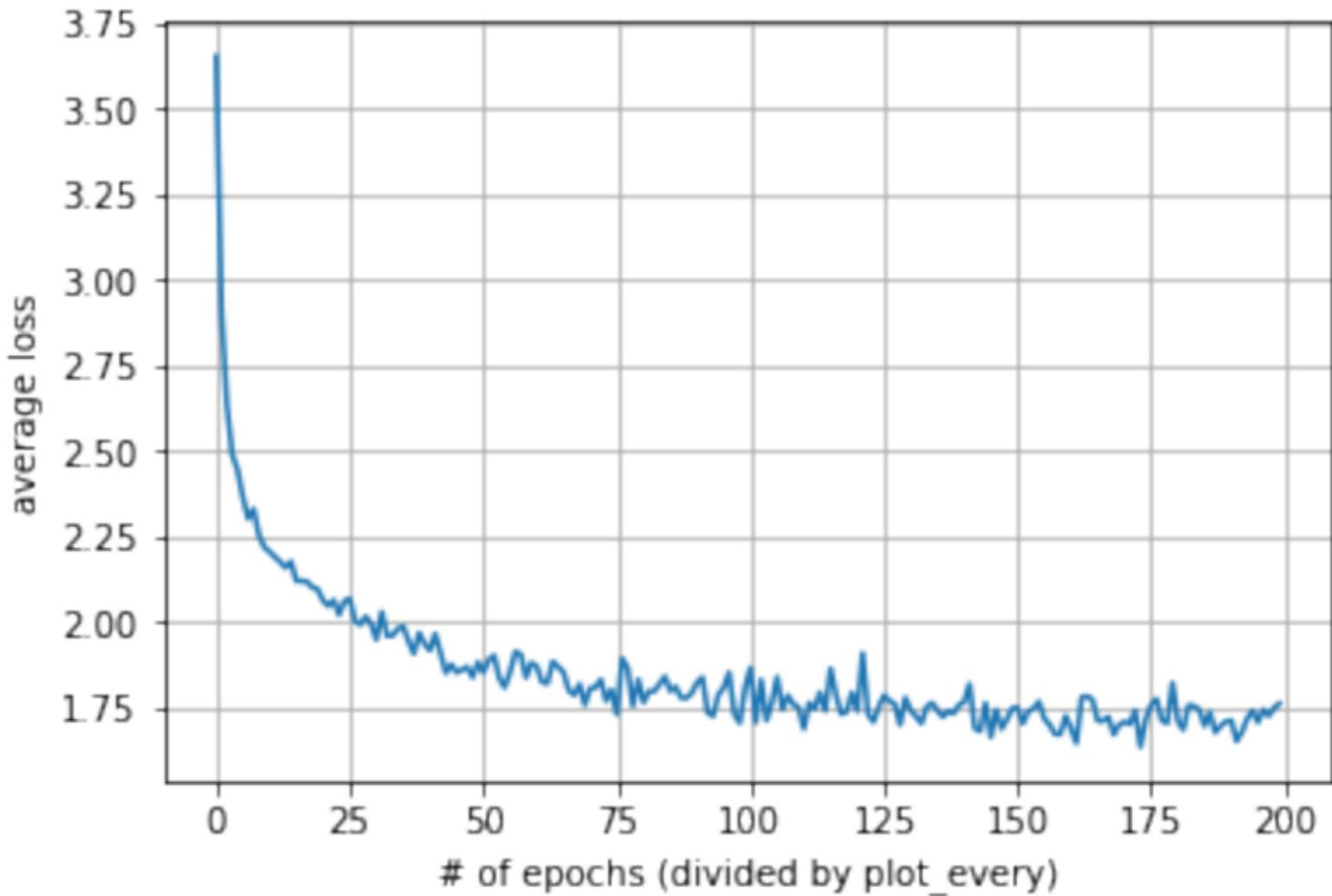
5.

b

3)

Generating Shakespeare Using a Character-level Language Model.

- One very noticeable advantage a char-based model would have over a word-based model is in terms of computational resources. The dimension $|V|$ is absolutely huge in comparison to the number of characters in the language [for most languages]. Working with matrices of a much smaller dimension leads to a very substantial improvement in terms of runtime & space required.
- Also, a char-based model would also be more flexible in its ability to treat language - anomalies [since the number of different characters is limited whilst that of words is possibly infinite].
- That same trait might also come as a disadvantage in some tasks where we wish to avoid gibberish-meaningless words from being generated.



4 PLEXITY:

• Denote: $q(s_i) = p(s_i | s_1, \dots, s_{i-1})$

$$\begin{aligned} & \left[-\frac{1}{m} \sum_{i=1}^m \log_b (q(s_i)) \right] = b \left[-\frac{1}{m} \log_b \left[\prod_{i=1}^m q(s_i) \right] \right] \\ & = b \log_b \left[\left[\prod_{i=1}^m q(s_i) \right]^{-\frac{1}{m}} \right] = \left[\prod_{i=1}^m q(s_i) \right]^{-\frac{1}{m}}. \end{aligned}$$

$b=2, e$ 1.778 027-1, b の $\log_b 1/2$ は

$$2 \left[-\frac{1}{m} \sum_{i=1}^m \log_2 (q(s_i)) \right] = e \left[-\frac{1}{m} \sum_{i=1}^m \ln (q(s_i)) \right]$$

①