

1 a)  $J_{\text{naive-soltnr}}(v_c, v, U) = -\log P(O=o | C=c)$   
 $= -\sum_{w \in \text{vocab}} y_w \log(\hat{y}_w) = -\log(\hat{y}_o)$

b)  $\hat{y} = P(O | C=c)$ , so  $\text{eqn. 2} = \text{eqn. 3}$

b) i)  $J_{\text{naive-soltnr}}(v_c, v, U) = -\log \frac{\exp(u_o^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)}$

$$\frac{dJ}{dv_c} = -\frac{1}{\hat{y}_o} \cdot \frac{\partial \hat{y}_o}{\partial v_c}$$

$$\frac{\partial \hat{y}_o}{\partial v_c} = \frac{\partial}{\partial v_c} \left( \frac{\exp(u_o^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} \right)$$

$$= \left( \sum_{w \in \text{vocab}} \exp(u_w^T v_c) \cdot \exp(u_o^T v_c) \cdot u_o \right)$$

$$= \left( \exp(u_o^T v_c) \cdot \sum_{w \in \text{vocab}} \exp(u_w^T v_c) \cdot u_w \right)$$

$$\left( \sum_{w \in \text{vocab}} \exp(u_w^T v_c) \right)^2$$

$$= \frac{\exp(u_o^T v_c) (u_o - u_w)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} = \text{softmax}$$

$$= \hat{y}_o u_o = \sum_{w \in \text{vocab}} \hat{y}_w u_w$$

$$\frac{dJ}{dv_c} = -\frac{1}{\hat{y}_o} \left( \hat{y}_o u_o - \sum_{w \in \text{vocab}} \hat{y}_w u_w \right)$$

$$= -u_o + \sum_{w \in \text{vocab}} \hat{y}_w u_w$$

$$= U(\hat{y} - y)$$



ii)  $\frac{dJ}{dv_c} = 0$  when  $\frac{1}{\gamma} = \gamma$   $\square$

iii) when  $\frac{dJ}{dv_c}$  subtracted from word  $v_c$ ,

it makes  $v_c$  closer to  $u_0$  (aka the correct class)  $\square$

c) When magnitude important, not lost vs  
when dir. " not lost,  $\square$

d)  $W=0: \frac{dJ}{du_0} = -\frac{1}{\gamma_0} \cdot \frac{d\gamma_0}{du_0}$

$$\frac{2\gamma_0}{2u_0} = \frac{2}{2u_0} \left( \frac{\exp(u_0^T v_c)}{\sum_{w \in \text{Vocab}} (u_w^T v_c)} \right)$$

$$= \frac{2}{2u_0} \left( \frac{\exp(u_0^T v_c)}{u_0^T v_c} \right)$$

$$= \frac{(u_0^T v_c) \cdot \exp(u_0^T v_c) - v_c \cdot \exp(u_0^T v_c)}{(u_0^T v_c)^2}$$

$$= \frac{\exp(u_0^T v_c) \cdot v_c ((u_0^T v_c) - 1)}{(u_0^T v_c)^2}$$

$$= \gamma_0 (v_c \cdot ((u_0^T v_c) - 1))$$

$$\frac{dJ}{du_0} = \boxed{-v_c (\gamma - 1)}$$

$$w \neq 0: \frac{\partial \hat{y}_c}{\partial u_w} = \frac{\partial}{\partial u_w} \left( \frac{\exp(u_0^T V_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T V_c)} \right)$$

$$= \frac{\exp(u_0^T V_c) \cdot 0 \cdot \sum_{w \in \text{vocab}} \exp(u_w^T V_c) - \exp(u_0^T V_c) \cdot \sum_{w \in \text{vocab}} \exp(u_w^T V_c) \cdot V_c}{\left( \sum_{w \in \text{vocab}} \exp(u_w^T V_c) \right)^2}$$

$$\frac{\partial \hat{y}_c}{\partial u_w} = -\frac{1}{\hat{y}_c} \cdot \hat{y}_c \cdot V_c$$

$$= \boxed{\hat{y}_w - V_c}$$

$$e) \frac{dJ}{dU} = \begin{bmatrix} \frac{\partial J(u_{c1}, U)}{\partial u_1} & \dots & \frac{\partial J(u_{cI}, U)}{\partial u_I} \end{bmatrix} \begin{bmatrix} \frac{\partial J(u_{c1}, U)}{\partial u_1} & \dots & \frac{\partial J(u_{cI}, U)}{\partial u_I} \end{bmatrix}$$

$$= \begin{bmatrix} \hat{y}_1 - V_c & \dots & \hat{y}_w - V_c \end{bmatrix}$$

2. a) i) Updates can only affect gradient by small amount  $\rightarrow$  less var  $\rightarrow$  less noise when training
- ii) adaptive lr  $\rightarrow$  bigger updates when grad mag.  $\uparrow$ , smaller if  $\downarrow$ ,  
move faster in right dir & slower in wrong dir.

b) i.)  $\mathbb{E}[\text{drop} \cdot [Y \cdot \theta]_i] = h_i$

$$= \gamma \mathbb{E}[d_i] \cdot h_i = h_i$$

$$\Rightarrow \gamma = \frac{1}{\mathbb{E}[d_i]} = \frac{1}{0 \cdot \text{pdrop} + 1 \cdot (1 - \text{pdrop})} = \frac{1}{1 - \text{pdrop}}$$

ii) Only during training, bc throwing away prod



3) a) (+3 given),  $n = 8$  words

Stack	Buffer	New dep.	Trans.
[root, presented, my.]	[findings, at, the, NLP,		SHIFT
[root, presented, my, findings]	conference]		SHIFT
[root, pres, findings]	[at, the, NLP, conf.]	findings $\rightarrow$ my	LEFT-ARC
[root, pres.]	"	pres. $\rightarrow$ findings	RIGHT-ARC
[root, pres., at]	[the, NLP, conf.]		SHIFT
[root, pres., at, the]	[NLP, conf.]		SHIFT
[root, pres., at, the, NLP]	[conf.]		SHIFT
[root, pres., at, the, NLP, conf.]			SHIFT
[root, pres., the, NLP, conf.]		conf. $\rightarrow$ at	LEFT
[root, pres., NLP, conf.]		conf. $\rightarrow$ the	LEFT
[root, pres., conf.]		conf. $\rightarrow$ NLP	LEFT
[root, pres.]		pres. $\rightarrow$ conf.	RIGHT
[root]		pres. $\rightarrow$ root	RIGHT

b)  $n = 8 \rightarrow 2n$  steps =  $8(2) = 16$  steps b/c  $n$  steps

to  $\rightarrow$  stack,  $n$  steps to remove.

c) Code

d) code  $\begin{matrix} 3 \times 4 \\ (3 \times 2) & (4 \times 4) & (4) \end{matrix} \quad 4 \times 2$

e) i)  $\frac{\partial h_i}{\partial x_j} = (xW + b)_i > 0 \times W.T$

ii)  $\frac{\partial CE(y, \hat{y})}{\partial \hat{y}_i} = \sum_j \frac{\partial CE(y, \hat{y})}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial \hat{y}_i} \quad (j \neq i)$

$= \frac{\partial CE(y, \hat{y})}{\partial \hat{y}_c} \frac{\partial \hat{y}_c}{\partial \hat{y}_i}$

$= \frac{y_i}{\hat{y}_c} \cdot \frac{\partial}{\partial \hat{y}_i}$

$= \frac{y_i}{\hat{y}_c} \cdot e^{l_i - l_j} \cdot 1 = \frac{e^{l_i}}{\sum_{j=1}^2 e^{l_j}} \cdot \frac{e^{l_i - l_j} (y_i)}{\hat{y}_c}$



iii) code

f) i) error: . . . verb error

incorr. dep.: . . . sitting →

corr. dep.: . . . sitting → blocked

ii) "

mod error

elements → most

crucial → most

iii) "

prep. phrase attach error

declined → decision

reasons → decision

iv) "

coord error

one → and

Quebec → and

g) P.D. tags can give signal as they often follow

certain words.