

1.a) i. "JORDAN is BEAUTIFUL"  
 PER OR LOC

"I TOOK A LOAN FROM GOLDMAN".  
 ORG OR PER

ii. IT MIGHT BE IMPORTANT TO USE THOSE WORDS TO IDENTIFY THE CONTEXT AND CLEAR THOSE AMBIGUITIES.

iii.  $\rightarrow$  THE FOLLOWING WORD<sup>pos</sup> (E.G., IF IT IS A VERB, THE CHANCES OF THE PRECEDING WORD BEING A NAMED ENTITY IS HIGHER).

$\rightarrow$  THE PREVIOUS WORD<sup>pos</sup> (E.G., IF IT IS A WORD SUCH AS "TO", "FROM", "IN", THE CHANCES OF THE FOLLOWING WORD BEING A NAMED ENTITY IS HIGHER).

1.b) i.  $e^{(t)}: 1 \times (2w+1) \times D$      $h^{(t)}: 1 \times H$

$W: (2w+1)D \times H$

$U: H \times C$

ii.  $e^{(t)} \cdot W \Rightarrow O((2w+1)HD)$   
 $h^{(t)} \cdot U \Rightarrow O(HC)$  } THOSE OPERATIONS DOMINATE THE COMPLEXITY

WE HAVE  $T$  WORDS IN THE SENTENCE.

THEREFORE, OVERALL COMPLEXITY IS  $O(TH \cdot (C + (2w+1)D))$

CS224N

1.d) i. BEST F1 SCORE: 0.83

TOKEN-LEVEL CONFUSION MATRIX:

ACTUAL \ PREDICTED	GO	PER	ORG	LOC	MISC	0
GO	2921	50	54	21	103	
PER	135	1657	99	55	146	
ORG	43	112	1866	19	54	
LOC	45	60	38	1011	114	
MISC	37	55	13	26	42628	
0						

THE MODEL OFTEN PREDICTS "0" WHEN IT'S ACTUALLY "ORG" OR "MISC", AND ALSO PREDICTS "PER" WHEN IT'S ACTUALLY "ORG" VERY FREQUENTLY

1.d)ii.

- If the window size is not large enough, the model might miss some long continuous named entities with connector words inside it.

Example:

$x$ :	May	15	v	Duke	of	Norfolk	...
$y^*$ :	0	0	0	ORG	ORG	ORG	...
$y'$ :	0	0	0	ORG	0	ORG	...

- The model might not have enough memory to understand context and discriminate between different types of named entities.

Example:

$x$ :	...	until	the	sixth	game,	when	Washington,	after	winning	...
$y^*$ :	...	0	0	0	0	0	PER	0	0	...
$y'$ :	...	0	0	0	0	0	LOC	0	0	...

2.a) i. window based model:  $(2w+1)DH$  (from  $W$ )

Not considering  $L$ , since it's the same for both

$H \times C$  (from  $U$ )

$H$  (from  $b_1$ )

$C$  (from  $b_2$ )

RNN:  $H \times H$  (from  $W_h$ )

$D \times H$  (from  $W_x$ )

$H$  (from  $b_1$ )

$H \times C$  (from  $U$ )

$C$  (from  $b_2$ )

So we have  $H^2 + DH + \cancel{H} + \cancel{HC} + \cancel{C} - \cancel{C} - \cancel{H} - \cancel{HC} - (2w+1)DH =$

$H \cdot (H - 2wD)$  more parameters.

2.a) iii. We need to go over each of the  $T$  words. So we will have a factor of  $O(T)$ . Now, for each word, the dominating operations are

$$\begin{matrix} h^{(t-1)} & & e^{(t)} & & h^{(t)} \\ 1 \times H & H \times H & 1 \times D & D \times H & L \times H & H \times C \end{matrix} \Rightarrow \begin{matrix} O(H^2) & O(DH) & O(HC) \end{matrix}$$

Total complexity:  $O(TH(H+D+C))$ .



2.6) i.

FROM	$y$	PER	0	PER	PER	PER	PER
	$\tilde{y}$	PER	0	0	0	0	0
TO	$y$	PER	0	PER	PER	PER	PER
	$\tilde{y}$	0	0	0	PER	PER	PER

THE ENTITY LEVEL  $F_1$  IS NOW 0 (DECREASED),  
 EVEN THOUGH THE CROSS ENTROPY COST  
 WOULD HAVE DECREASED (ASSUMING  $\hat{y}_{\text{PER}} = 1$  IF  $\tilde{y} = \text{PER}$ ,  $\hat{y}_0 = 1$  IF  $\tilde{y} = 0$ ).

2.6) ii. IT'S DIFFICULT TO DIRECTLY OPTIMIZE FOR  $F_1$  BECAUSE  
 TO IT IS A NON-DIFFERENTIABLE FUNCTION (NEITHER PRECISION OR RECALL  
 ARE).

2. d) i.

Since those new  $y$ s are filled with 0s, and

$$CE(y^{(t)}, \hat{y}^{(t)}) = \sum_i y_i^{(t)} \log(\hat{y}_i^{(t)}) = \sum_i 0 \cdot \log(\hat{y}_i^{(t)}) = 0,$$

the loss  $J$  doesn't change. However, the gradients do change: as we saw before, if  $\hat{y} = \text{softmax}(\theta)$  and

$$J = CE \text{ loss}, \quad \frac{\partial J}{\partial \theta} = (\hat{y} - y).$$

Even though  $y$  is filled with 0s,  $\hat{y}$  isn't (since it is the output of a softmax).

So if we don't apply the mask, we will see some non-zero gradient updates on the parameters when we do

backpropagation after the pass of one of those "new" pair of vectors  $(x, y)$ . Masking solves the problem by having on  $m^{(t)} = \mathbb{1} - \mathbb{1}\{t > T\}$  multiplying the loss (and, thus, the gradients) for each pass  $t = 1, \dots, T$ . That would set to 0 both the loss and the gradients for those "new" pairs  $(x, y)$ , solving the issue.

2.g) i, ii. • The model only reads the words until the current word  $t$ , i.e., just learn the past context.

Example:

$x$	...	the	Stars	and	Stripes	flag	was	...
$y^*$	...	0	MISC	MISC	MISC	0	0	...
$y'$	...	0	LOC	0	MISC	0	0	...

Model extension: build a bidirectional recurrent neural network.

• For long sentences, the model might be losing information about the context as it reaches the end of the sentence, which could be due vanishing gradients or just the fact that it might not be capturing the relevant information that it should save for the future.

Example:

$x$	... [long sentence] ...	government	of	Fulgencio	Batista.
$y^*$	...	0	0	PER	PER 0
$y'$	...	0	0	ORG	ORG 0

Model extensions: using ReLU as the activation function, or even move to more sophisticated models such as GRUs or LSTMs.

3. a) possible cases:

→ First element is 0:  $h^1 = 0 = \sigma(0 + 0 + b_n) \Rightarrow b_n \leq 0$

→ First element is 1:  $h^1 = 1 = \sigma(u_n + 0 + b_n) \Rightarrow u_n + b_n > 0$

→ Keep previous element 1:  $\begin{cases} h^t = 1 = \sigma(u_n + w_n + b_n) \Rightarrow w_n + b_n > 0 \\ h^t = 1 = \sigma(u_n + w_n + b_n) \Rightarrow u_n + w_n + b_n > 0 \end{cases}$

→ Keep previous element 0 (analogous to 1st case).

Setting values:  $b_n = -1$ ;  $u_n = 2$ ;  $w_n = 2$

ii.  $h^1 = 0$ ;  $h^2 = 1$ ; possible cases:

→  $h^1 = 0$ ;  $h^2 = 1$  =



# CS 224N

3. a) ii.  $w_n = u_n = b_n = b_z = b_h = 0$ :

$$z^t = \sigma(x^t u_z + h^{t-1} w_z)$$

$$x^t = \sigma(0 + 0 + 0) = \sigma(0) = 0$$

$$\tilde{h}^t = \tanh(x^t u_h + 0 + 0)$$

$$h^t = z^t \cdot h^{t-1} + (1 - z^t) \cdot \tilde{h}^t$$

Possible cases:

→ Keep previous 0:  $h^t = 0 = \underbrace{\sigma(0)}_{z^t} \cdot 0 + (1 - 0) \cdot \tanh(0) = 0 \Rightarrow \text{o.k.}$

→  $x^t = 0, h^{t-1} = 1$ :  $h^t = 1 = \sigma(w_z) \cdot 1 + (1 - \sigma(w_z)) \cdot \overbrace{\tanh(0)}^{=0} \Rightarrow w_z > 0$

→  $x^t = 1, h^{t-1} = 0$ :  $h^t = 1 = \sigma(u_z) \cdot 0 + (1 - \sigma(u_z)) \cdot \tanh(u_h) \Rightarrow \begin{cases} u_z \leq 0 \\ u_h > 0 \end{cases}$

Possible solution:

$$w_z = 1; u_z = -1; u_h = 1; w_h = 0$$

$w_h \Rightarrow \text{any value}$

3. b)i.

possible cases:

$$\rightarrow h^{t-1}=0, x=0: h^t=0 = \sigma(0+0+b_n) \Rightarrow b_n \leq 0 \quad \textcircled{I}$$

$$\rightarrow h^{t-1}=0, x=1: h^t=1 = \sigma(u_n+0+b_n) \Rightarrow u_n+b_n > 0 \quad \textcircled{II}$$

$$\rightarrow h^{t-1}=1, x=0: h^t=1 = \sigma(0+w_n+b_n) \Rightarrow w_n+b_n > 0 \quad \textcircled{III}$$

$$\rightarrow h^{t-1}=1, x=1: h^t=0 = \sigma(u_n+w_n+b_n) \Rightarrow u_n+w_n+b_n \leq 0 \quad \textcircled{IV}$$

$$-\textcircled{IV} + \textcircled{II} - \textcircled{III}: -u_n - w_n - b_n + u_n + b_n + w_n + b_n \geq 0 \Rightarrow b_n \geq 0.$$

From  $\textcircled{I}$ , we get that  $b_n = 0$ .

No, looking again at  $\textcircled{II} + \textcircled{III}$ , we get  $u_n + w_n > 0$ , which is impossible since, from  $\textcircled{IV}$ ,  $u_n + w_n \leq 0$ .

CS 224N

3.b)ii.  $w_n = u_n = b_z = b_n = 0$

$$z^{(t)} = \sigma(x^{(t)} u_z + h^{(t-1)} w_z)$$

$$x^{(t)} = \sigma(b_n)$$

$$\tilde{h}^{(t)} = \tanh(x^{(t)} u_n + \sigma(b_n) w_n \cdot h^{(t-1)})$$

$$h^{(t)} = z^{(t)} h^{(t-1)} + (1 - z^{(t)}) \tilde{h}^{(t)}$$

Case:

$$\Rightarrow h^{(t-1)} = 0; x^{(t)} = 0 \Rightarrow h^{(t)} = 0 = \sigma(0) \cdot 0 + \underbrace{1}_{(1-\sigma(0))} \cdot \underbrace{0}_{\tanh(0)} = 0$$

$$\Rightarrow h^{(t-1)} = 0; x^{(t)} = 1:$$

$$h^{(t)} = 1 = \underbrace{0}_{\sigma(u_z) \cdot 0} + (1 - \sigma(u_z)) \cdot \tanh(u_n + 0) \Rightarrow u_z \leq 0, u_n > 0$$

$$\Rightarrow h^{(t-1)} = 1; x^{(t)} = 0:$$

$$h^{(t)} = 1 = \sigma(w_z) + (1 - \sigma(w_z)) \cdot \tanh(\sigma(b_n) w_n) \Rightarrow \begin{cases} w_z > 0 \\ \text{or} \\ b_n > 0, w_n > 0 \end{cases}$$

$$\Rightarrow h^{(t-1)} = 1; x^{(t)} = 1:$$

$$h^{(t)} = 0 = \sigma(u_z + w_z) + (1 - \sigma(u_z + w_z)) \cdot \tanh(u_n + \sigma(b_n) w_n) \Rightarrow \begin{cases} u_z + w_z \leq 0 \\ \text{and} \\ u_n + \sigma(b_n) w_n \leq 0 \end{cases}$$

POSSIBLE SOLUTION:

$\underline{u_n} = 1$	$\underline{u_z} = -1$	$\underline{u_n} = 1$
$\underline{w_z} = 1$	$\underline{w_n} = -1$	