

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. → 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).

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

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

$V: H \times C$

ii.  $e^{(t)}.W \Rightarrow O((2w+1)HD)$        $h^{(t)}.V \Rightarrow O(HC)$       } 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

Predicted		GU	PER	ORG	LOC	Misc	O
Go		PER	ORG	LOC	Misc	O	
PER	2921	50	54	21	103		
ORG	135	1657	99	55	146		
LOC	43	112	1866	19	54		
MISC	45	60	38	1011	114		
O	37	55	13	26	42628		

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

# CS 224N

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	...
---------	---	---	---	-----	-----	-----	-----

$y^{\prime}$ :	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	...
---------	-----	---	---	---	---	---	-----	---	---	-----

$y^{\prime}$ :	...	0	0	0	0	0	LOC	0	0	...
----------------	-----	---	---	---	---	---	-----	---	---	-----

2.a) i. Window based model:  $(2w+L)DH$  (from  $W$ )  
 Not considering  $L$ , since it's the same for both RNN:  $H \times H$  (from  $W_h$ )

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

$H$  (from  $b_1$ )

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

$C$  (from  $b_2$ )

No we have  $H^2 + DH + H + HC + C - C - H - HC - (2w+L)DH = H \cdot (H - 2wD)$  more parameters.

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

$$\begin{array}{lll} h^{(t-1)} W_h \Rightarrow O(H^2) & e^{(t)} \cdot W_x \Rightarrow O(DH) & h^{(t)} U \Rightarrow O(HC) \\ 1 \times H \quad H \times H & 1 \times D \quad D \times H & L \times H \quad H \times C \end{array}$$

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

2.b)i.

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

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

2.b)ii. IT'S DIFFICULT TO DIRECTLY OPTIMIZE FOR  $F_1$  BECAUSE IT IS A NON-DIFFERENTIABLE FUNCTION (NEITHER PRECISION OR RECALL ARE CATEGORICAL, SO THERE IS NO GRADIENT OF THE LOSS FUNCTION).

( $y_0, y_{\text{PER}}$  ALLOW GRADIENTS)

2.d) i.

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

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

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

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

an  $m^{(t)} = 1 - \mathbb{1}\{t > T\}$  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.

# CS 224N

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.

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

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

3. a) Possible cores:

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

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

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

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

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

Difficult to implement; Realistic?

$\rightarrow h^t = \sigma(u_n + w_n + b_n) = \sigma(2 + 2 - 1) = 1$

# CS 224N

3.a) ii.  $w_n = u_n = b_n = b_z = b_w = \emptyset$ :

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

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

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

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

Possible cases:

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

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

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

Possible solution:

$w_z = 1; \quad u_z = -1; \quad u_n = 1; \quad w_n = \emptyset$
---

$w_n \Rightarrow \text{any value}$

3. b)i.

Possible cases:

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

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

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

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

$$-\textcircled{IV} + \textcircled{II} - \textcircled{III}: -\cancel{\mu_n} - \cancel{w_n} - \cancel{b_n} + \cancel{\mu_n} + \cancel{b_n} + \cancel{w_n} + b_n > 0 \Rightarrow b_n > 0.$$

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

Now, looking again at \textcircled{II} + \textcircled{III}, we get  $\mu_n + w_n > 0$ , which is impossible since, from \textcircled{IV},  $\mu_n + w_n < 0$ .

# CS 224N

3.b) ii.  $w_n = u_n = b_3 = b_n = 0$

$$z^{(t)} = \sigma(u^{(t)} u_3 + h^{(t-1)} w_3)$$

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

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

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

Case 1:

$$\Rightarrow h^{(t-1)} = 0; u^{(t)} = 0 \Rightarrow h^{(t)} = \sigma(0) \cdot 0 + \overbrace{\frac{1}{(1-\sigma(0))}}^0 \cdot \overbrace{\tanh(0)}^0 = 0$$

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

$$h^{(t)} = 1 = \sigma(u_3) \cdot 0 + (1 - \sigma(u_3)) \cdot \tanh(u_n + 0) \Rightarrow u_3 \leq 0, u_n > 0$$

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

$$h^{(t)} = 1 = \sigma(w_3) \dots + (1 - \sigma(w_3)) \cdot \tanh(\sigma(b_n) w_n) \Rightarrow \begin{cases} w_3 > 0 \\ \text{or} \\ b_n > 0, w_n > 0 \end{cases}$$

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

$$h^{(t)} = 1 = \sigma(u_3 + w_3) + (1 - \sigma(u_3 + w_3)) \cdot \tanh(u_n + \sigma(b_n) w_n) \Rightarrow \begin{cases} u_3 + w_3 \leq 0 \\ \text{and} \\ u_n + \sigma(b_n) w_n \leq 0 \end{cases}$$

Possible Solution:

$b_3 = 1$	$u_3 = -1$	$u_n = 1$
$w_3 = 1$	$w_n = -1$	

2.d) ii.

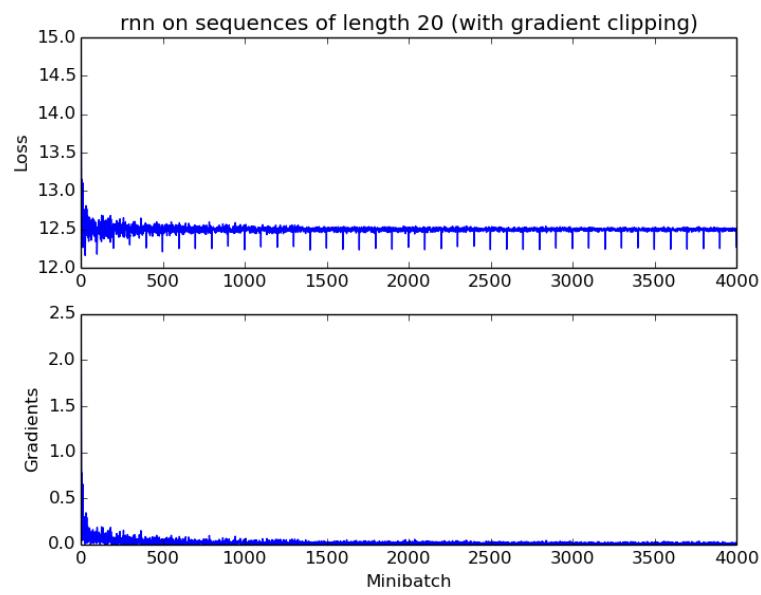
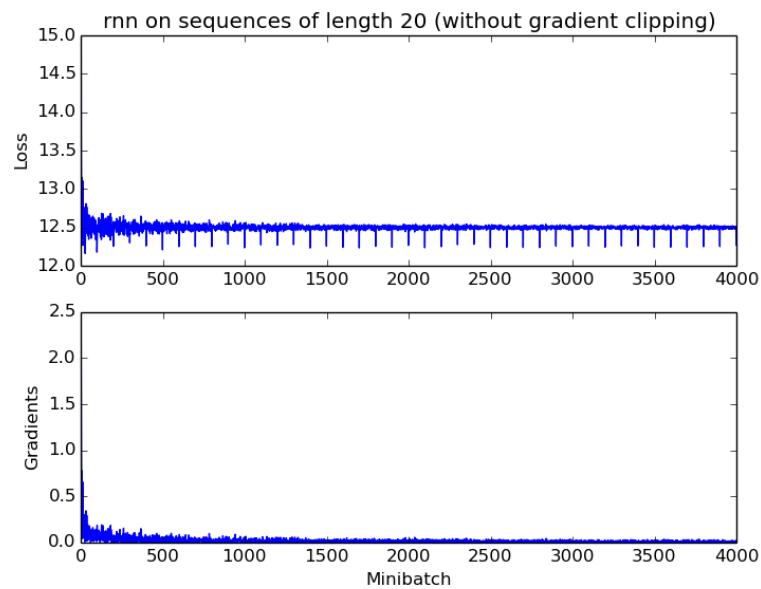
If we did not use masking, the loss would increase, since we would be adding M-T terms of the form  $\hat{y}^{(t)}$ ,  $t > T$ .

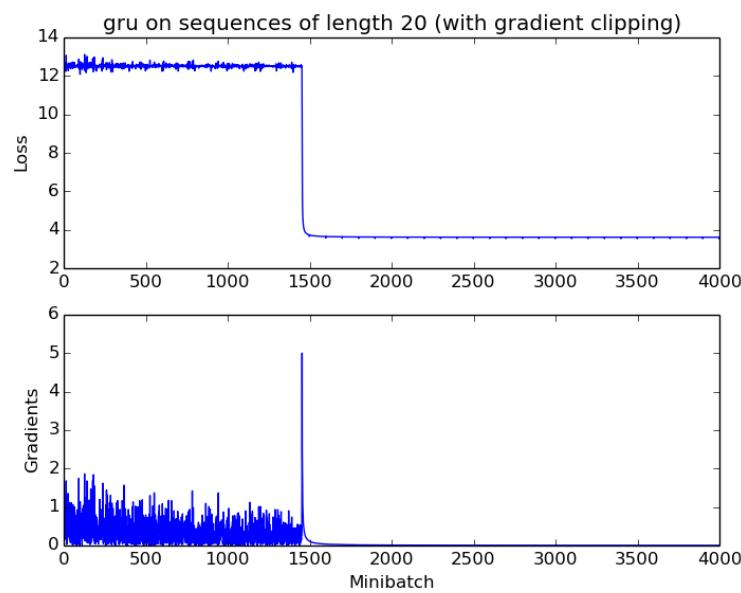
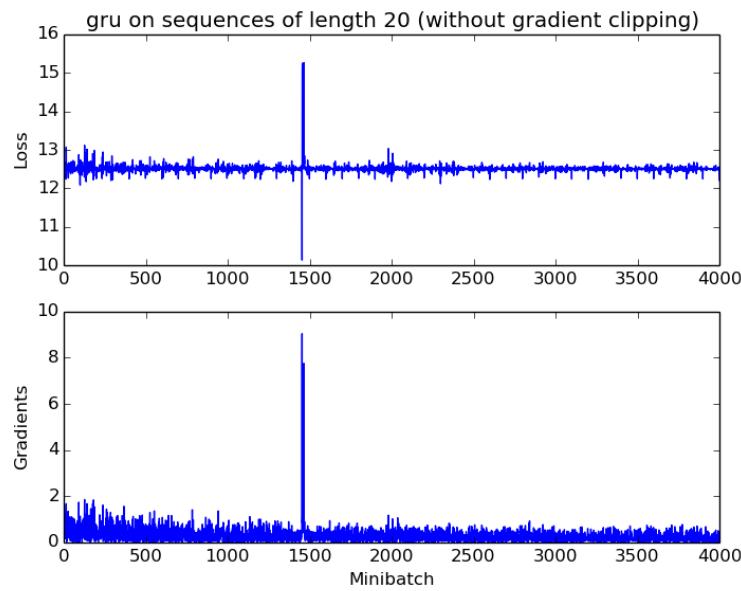
Also, the gradient would change. As we saw before, if

$\hat{y} = \text{softmax}(\theta)$  and  $J = \text{CE loss}$ ,  $\frac{\partial J}{\partial \theta} = (\hat{y} - y)$ . So not only  $y$  has a 1 in the "null" position, but  $\hat{y}$  is filled with nonzero elements as well. This first error  $\hat{y}$  backpropagates and will end up affecting most parameters' gradients.

Masking solves the problem by having an  $m^{(t)} = \mathbb{I}\{t \leq T\}$  multiplying each term of the loss (and, thus, the gradients), for each pass  $t = 1, \dots, M$ . That would set to 0 both loss and gradients for  $t = T+1, \dots, M$  for  $(x^{(t)}, y^{(t)})$ , solving the issue (i.e., the new gradients and loss would be equivalent w.r.t. if we didn't do the augmentation).

### Question 3.d)





### **Question 3.e)i.**

It seems like the RNN model experience vanishing gradients, since the gradients become very small and the model get stuck in some local minima. On the other hand, the GRU model seems to experience exploding gradient, since at some point it reaches a loss significantly better than before, but a large gradient pushes it to a worse local minimia region. Clipping does help, as we can see from the last pair of graphs: after it reaches this better loss, the clipped gradient is not large enough to push it to the worse local minima, and then we stay in the same region, eventually converging to a significantly better loss.

### **Question 3.e)ii.**

GRU does better since it is able to achieve a much lower loss. This doesn't necessarily guarantees a better generalization error, but the fit is definitely better, we achieve a much lower error and the model seems to allow for better accuracy.