

Student Name: Chaoran Li      Student NetID: cxl190012

University of Texas at Dallas  
Department of Computer Science  
CS6320 – Natural Language Processing  
Spring 2021

Instructor: Dr. Sanda Harabagiu

### Take-Home Mid-Term Exam

Issued: March 25<sup>th</sup> 2021

Due: March 27<sup>th</sup> 2021 – before MidNight

**Submit in eLearning as PDF file**

**DO NOT DELETE ANYNITHING, Simply add your answers!!!!**

If you submit only the solution with no problems, you will receive 0 points!!!!

---

**Instructions:** Do not communicate with anyone in any shape or form. This is an independent exam. Do not delete any problem formulation, just attach your answer in the space provided. If the problem is deleted and you send only the answer, you shall receive ZERO points. If you do not write your name and netid, it will be considered that you did not submit your Midterm exam, and will obtain ZERO points.

Copy and paste the Midterm Exam into a Word document, enter your answers (either by typing in Word, or by inserting a VERY CLEAR picture of your hand-written solution) and transform the file of the exam into a PDF format. If we cannot clearly read the picture, you will get ZERO for that answer! If you create an enormous file for your final exam (i.e. larger than 5 Mbytes) you will receive ZERO for your entire exam. Please follow the instructions from the attached **instructions\_submission\_EXAM.pdf** file to make sure your final pdf file is of reasonable size.

If you will use a pencil instead of a black pen, you will receive a ZERO for the entire final exam.

Make sure that you insert EACH answer immediately after EACH question. Failure to do so will result in ZERO points for the entire exam! Submit the PDF file with the name **Final\_Exam\_netID.pdf**, where netID is your unique netid provided by UTD. If you submit your exam in any other format your will receive ZERO points.

The MidTerm exam shall be submitted in eLearning before the deadline. No late submissions shall be graded! Any cheating attempt will determine the ENTIRE grade of the final exam to become ZERO.

Write your answers immediately after the problem statements. If you enter multiple possible answers to the same problem –the most incorrect answer shall be selected and graded (as you are not sure about which answer is the correct one!)

### Problem 1 Language Models [ TOTAL: 50 points ]

You are given the Bigram probability matrix **M** containing evaluations on the Training corpus of the maximum likelihood estimations of the probabilities of the bigrams for the words "Tom", "talks", "a", "lot" and "sometimes". The corpus has 2500 sentences (with not punctuation signs!) and a vocabulary  $|V|=25$ . You also know that the unigram counts for the words are: Tom = 1500, a = 1000, talks = 2000; sometimes=550 and lot = 400. The matrix **M** is:

	<s>	Tom	talks	a	lot	sometimes	</s>
<s>	0.0	0.65	0.0	0.15	0.0	0.2	0.0
Tom	0.0	0.0	0.7	0.0	0.0	0.3	0.0
talks	0.0	0.0	0.0	0.65	0.0	0.15	0.2
a	0.0	0.0	0.0	0.0	0.9	0.1	0.0
lot	0.0	0.2	0.1	0.0	0.0	0.6	0.1
sometimes	0.0	0.2	0.5	0.0	0.0	0.0	0.3
</s>	1	0.0	0.0	0.0	0.0	0.0	0.0

- a) [ 10 points ] Reconstruct the bigram counts matrix from the bigram probabilities.  
Show how you obtained your calculations!

Problem 1 a) : Solution:

MLE n-grams (without smoothing) :

$$P(w_n | w_{n-n+1}^{n-1}) = C(w_{n-n+1}^{n-1}, w_n) / C(w_{n-n+1}^{n-1})$$

$$\Rightarrow C(w_{n-1}, w_n) = P(w_{n-1}, w_n) C(w_{n-1})$$

Besides, 2500 sentences => 2500 <s> and 2500 </s>

	<s>	Tom	talks	a	lot	sometimes	</s>
<s>	0.0x1500	0.65x2500	0.0x2500	0.15x2500	0.0x2500	0.2x2500	0.0x2500
Tom	0.0x1500	0.0x1500	0.7x1500	0.0x1500	0.0x1500	0.3x1500	0.0x1500
talks	0.0x2000	0.0x2000	0.0x2000	0.65x2000	0.0x2000	0.15x2000	0.2x2000
a	0.0x1000	0.0x1000	0.0x1000	0.0x1000	0.9x1000	0.1x1000	0.0x1000
lot	0.0x400	0.2x400	0.1x400	0.0x400	0.0x400	0.6x400	0.1x400
sometimes	0.0x550	0.2x550	0.5x550	0.0x550	0.0x550	0.0x550	0.3x550
</s>	1x2500	0.0x1500	0.0x1500	0.0x2500	0.0x2500	0.0x2500	0.0x2500

Hence, we can get the below results:

	<b>&lt; s &gt;</b>	<b>Tom</b>	<b>talks</b>	<b>a</b>	<b>lot</b>	<b>sometimes</b>	<b>&lt; /s &gt;</b>
<b>&lt; s &gt;</b>	0	1625	0	375	0	500	0
<b>Tom</b>	0	0	1050	0	0	450	0
<b>talks</b>	0	0	0	1300	0	300	400
<b>a</b>	0	0	0	0	900	100	0
<b>lot</b>	0	80	40	0	0	240	40
<b>sometimes</b>	0	110	275	0	0	0	165
<b>&lt; /s &gt;</b>	2500	0	0	0	0	0	0

- b) [5 points] Using the bigram probability matrix, compute the probability of the sentence “*Tom talks a lot sometimes*” and compare it to the probability of the sentence “*Sometimes Tom talks a lot*”.

Problem 1 b) : Solution:

S1: Tom talks a lot sometimes.

$$P_1 = P(\text{Tom} | \text{< s >}) P(\text{talks} | \text{Tom}) P(\text{a} | \text{talks}) P(\text{lot} | \text{a}) P(\text{sometimes} | \text{lot}) P(\text{< /s >} | \text{sometimes}) \\ = 0.65 \times 0.7 \times 0.65 \times 0.9 \times 0.6 \times 0.3 = 0.0479115$$

S2: Sometimes Tom talks a lot.

$$P_2 = P(\text{sometimes} | \text{< s >}) P(\text{Tom} | \text{sometimes}) P(\text{talks} | \text{Tom}) P(\text{a} | \text{talks}) P(\text{lot} | \text{a}) P(\text{< /s >} | \text{lot}) \\ = 0.2 \times 0.2 \times 0.7 \times 0.65 \times 0.9 \times 0.1 = 0.001638$$

“Tom talks a lot sometimes” has a higher possibility.

- c) [10 points] You also know that in the training corpus there is only one other word in the vocabulary that forms bigrams with  $w_1=\text{Tom}$ ,  $w_2=\text{talks}$ ,  $w_3=a$ ,  $w_4=\text{lot}$ , namely  $w_5=\text{home}$ .

The unigram count for  $w_5$  is 15000 in the Training corpus. You also know the following bigram probabilities from the Training corpus:

$\langle s \rangle$	Tom	talks	a	lot	home	$\langle /s \rangle$
0.0						
$\langle s \rangle$	0.7	0.0	0.1	0.1	0.1	0.0
Tom	0.0	0.0	0.75	0.1	0.1	0.05
talks	0.0	0.0	0.0	0.65	0.2	0.1
a	0.0	0.0	0.0	0.0	0.3	0.6
lot	0.0	0.2	0.4	0.3	0.0	0.05
home	0.0	0.5	0.0	0.1	0.3	0.0
$\langle /s \rangle$	1.0	0.0	0.0	0.0	0.0	0.0

Considering the following test set:

TEST: [Tom talks home a lot. Home Tom talks. Home talks.]

Show how you compute the perplexity of this bigram model smoothed with the Laplace method (**5 points**) and compare it to the perplexity of the language model when no smoothing was applied (**5 points**)

Problem 1 (c): Solution:

For Bigram,  $PP(w) = P(w_1, \dots, w_N)^{\frac{1}{N}} = \left[ \prod_{i=1}^N \frac{1}{P(w_i | w_{i-1})} \right]^{\frac{1}{N}}$

TEST:  $\langle s \rangle$  Tom talks home a lot  $\langle /s \rangle$   $\langle s \rangle$  Home Tom talks  $\langle /s \rangle$

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6 \quad w_7 \quad w_8 \quad w_9 \quad w_{10} \quad w_{11} \quad w_{12}$

$\langle s \rangle$  Home talks  $\langle /s \rangle$

$w_{13} \quad w_{14} \quad w_{15} \quad w_{16}$

$\therefore PP(w) = \left[ \left( \frac{1}{P(\text{Tom} | \langle s \rangle)} \cdot \frac{1}{P(\text{talks} | \text{Tom})} \cdot \frac{1}{P(\text{home} | \text{talks})} \cdot \frac{1}{P(\text{a} | \text{home})} \cdot \frac{1}{P(\text{lot} | \text{a})} \cdot \frac{1}{P(\langle /s \rangle | \text{lot})} \right) \cdot \frac{1}{P(\langle s \rangle | \langle s \rangle)} \cdot \frac{1}{P(\text{home} | \langle s \rangle)} \cdot \frac{1}{P(\text{Tom} | \text{home})} \cdot \frac{1}{P(\text{talks} | \text{Tom})} \cdot \frac{1}{P(\langle /s \rangle | \text{talks})} \cdot \frac{1}{P(\langle s \rangle | \langle /s \rangle)} \right]^{\frac{1}{20}}$

~~$P(\langle s \rangle | \langle s \rangle)$~~   ~~$P(\text{home} | \langle s \rangle)$~~   ~~$P(\text{Tom} | \text{home})$~~   ~~$P(\text{talks} | \text{Tom})$~~   ~~$P(\langle /s \rangle | \text{talks})$~~   ~~$P(\langle s \rangle | \langle /s \rangle)$~~

We can get all no smoothing probabilities from given table. But we need to calculate Laplace smoothing probabilities. Fortunately, we can merge some replicated items. Besides,  $|V| = 25$

$P_{\text{Laplace}}(w_i) = \frac{(i+1)}{N+V}, \quad P_{\text{MLE}}(w_i) = \frac{(i)}{N}, \quad (\text{Use } P_L \text{ short for Laplace and } P_{\text{MLE}} \text{ short for } P_{\text{MLE}} \text{ below})$

$\therefore P_L(w_i) = \frac{P_{\text{MLE}}(w_i)N+1}{N+V}$

$$\begin{aligned}
 P_L(Tom|<s>) &= \frac{0.7 \times 2500 + 1}{2500 + 25} = 0.693465 & P_L(talks|Tom) &= \frac{0.75 \times 1500 + 1}{1500 + 25} = 0.738361 \\
 P_L(home|talks) &= \frac{0.1 \times 12000 + 1}{12000 + 25} = 0.099259 & P_L(a|home) &= \frac{0.1 \times 1500 + 1}{1500 + 25} = 0.099900 \\
 P_L(lot|a) &= \frac{0.3 \times 1000 + 1}{1000 + 25} = 0.293658 & P_L(</s>|lot) &= \frac{0.05 \times 400 + 1}{400 + 25} = 0.049412 \\
 P_L(<s>|</s>) &= \frac{1.0 \times 2500 + 1}{2500 + 25} = 0.990495 & P_L(home|<s>) &= \frac{0.1 \times 2500 + 1}{2500 + 25} = 0.099406 \\
 P_L(Tom|home) &= \frac{0.5 \times 1500 + 1}{1500 + 25} = 0.499234 & P_L(talks|Tom) &= \frac{0.75 \times 1500 + 1}{1500 + 25} = 0.738361 \\
 P_L(</s>|talks) &= \frac{0.05 \times 2000 + 1}{2000 + 25} = 0.049876 & P_L(<s>|</s>) &= 0.990495 \\
 P_L(home|<s>) &= 0.099406 & P_L(talks|home) &= \frac{0 \times 1500 + 1}{1500 + 25} = 0.000066 \\
 P_L(\cancel{</s>}|talks) &= P(</s>|talks) = 0.049876
 \end{aligned}$$

$\therefore PP_{\text{Laplace}}(w) = [0.693465 \times 0.738361 \times 0.099259 \times 0.999000 \times 0.293658 \times 0.049412 \times 0.990495]$

$\approx [0.099406 \times 0.499234 \times 0.738361 \times 0.049876 \times 0.990495 \times 0.099406 \times 0.000066 \times 0.049876]^{1/8} = 6.843318$

Since  $P_M(\text{talks}|home) = 0$

[connect to Problem 1 c)]

$$PP_{\text{no-smoothing}}(w) = [\frac{1}{6}]^{1/8} = \infty$$

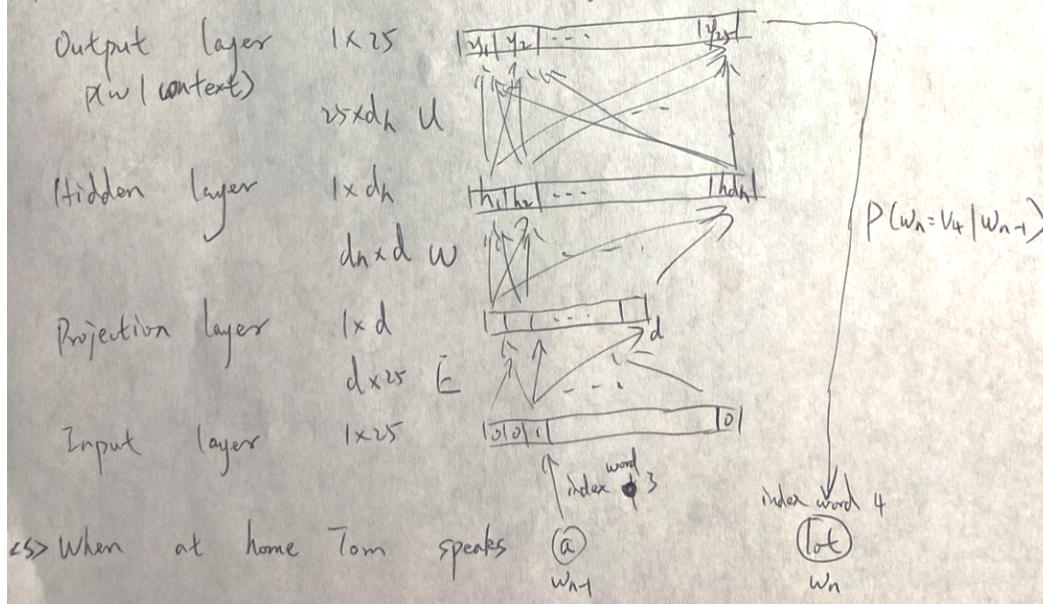
$$\therefore PP_{\text{Laplace}}(w) = 6.843318$$

$PP_{\text{no-smoothing}}(w) = \infty$   
 That is why we apply Laplace smoothing for calculating perplexity

We have  $PP_{\text{Laplace}}(w) = 6.843318$  and  $PP_{\text{no-smoothing}}(w) = \infty$ .

- d) [10 points] Draw a non-recursive neural architecture that allows you to (1) learn a neural language model from the training set and (2) learn in the same time the embeddings of the words. Exemplify the neural architecture working on the word sequence "When at home Tom speaks a" and show how it will predict the word "lot". [3 points] Detail the parameters of the neural model and explain how they are learned. Write the equations that define the working of this neural architecture [7 points]

Problem 1 d): Solution:  $|V|=25$  bigram:  $N=2-1=1$



Equations:

$e_i = \tilde{e}_i$	$\left. \begin{array}{l} \text{Connect to Problem 1 d)} \\ \text{Equations} \end{array} \right\}$
$h = \sigma(W_e h + b)$	
$z = Uh$	
$y = \text{softmax}(z)$	

- e) [15 points] Draw a recursive neural architecture that uses two stacked bi-directional GRU layers for learning a neural language model that operates on the word sequence "Sometimes at home Tom speaks a lot". [3 points] Detail the parameters of the neural model and explain how teacher forcing can be used and why. [5 points] Write the equations that define the working of this neural architecture [7 points]

Problem 1 e: Solution:

Output layer  $1 \times 7$

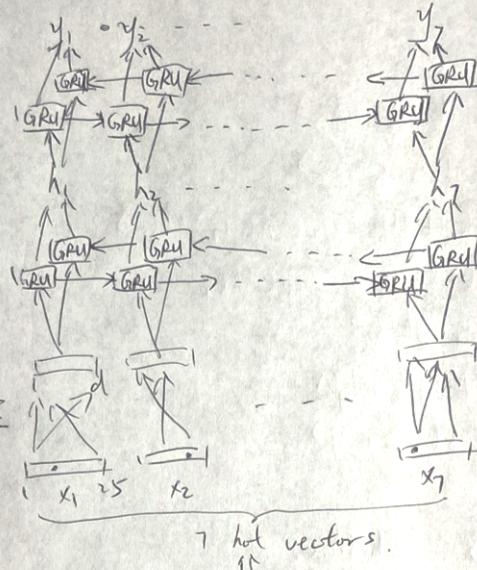
Bi-GRU layer

(Hidden layer  $\times$ )

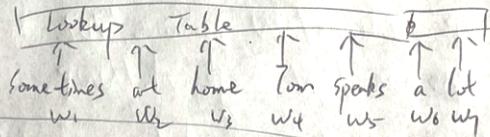
Bi-GRU layer

Projection layer  $1 \times 7 \times d$

Input layer  $1 \times 25 \times 7$



Words:  $|w|=7$



Teacher Forcing:

$L_{CE}(\hat{y}^t, y^t) = -\log \hat{y}_{w_{t+1}}^t$  The correct distribution  $y$  comes from knowing the next word.

At time  $t$ , the CE loss is the negative log probability assigned to the next word in the training sequence.

Equations:

[connects to Question 1 e]

$$x_t = p_c \times \text{Wembedding}$$

$$r_{t+1} = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z)$$

$$\tilde{h}_t = \tanh(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$
~~$$h_t = W^t h_m + w^b h_o^b$$~~
~~$$\odot h_t = \sigma(w^t e_t + b_t)$$~~

$$r_{t+2} = \sigma(W_r h_{t+2} + U_r h_{t+1} + b_r)$$

$$z_{t+2} = \sigma(W_z h_{t+2} + U_z h_{t+1} + b_z)$$

$$\tilde{h}_{t+2} = \tanh(W_h h_t + U_h (r_{t+2} \odot h_{t+1}) + b_h)$$

$$h_{t+2} = (1 - z_{t+2}) \odot h_{t+1} + z_{t+2} \odot \tilde{h}_{t+2}$$

$$e_t = W^t h_m + w^b h_o^b$$

$$P(y=j|e) = \frac{\exp(w_j^T e + b_j)}{\sum_i \exp(w_i^T e + b_i)}$$

### **Problem 2 (POS Tagging) (30 points)**

- a) (5 points) Using the Penn Treebank Part-of-Speech (POS) tag-set, manually assign tags for the following sentences:

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &amp;</i>
CD	cardinal number	<i>one, two, three</i>	TO	“to”	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential ‘there’	<i>there</i>	VB	verb, base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb, past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb, gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VBN	verb, past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb, non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb, 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, singular	<i>IBM</i>	\$	dollar sign	<i>\$</i>
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	“	left quote	<i>‘ or “</i>
POS	possessive ending	<i>'s</i>	”	right quote	<i>’ or ”</i>
PRP	personal pronoun	<i>I, you, he</i>	(	left parenthesis	<i>[, (, {, &lt;</i>
PRP\$	possessive pronoun	<i>your, one's</i>	)	right parenthesis	<i>], ), }, &gt;</i>
RB	adverb	<i>quickly, never</i>	,	comma	<i>,</i>
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	<i>. ! ?</i>
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	<i>: ; ... - -</i>
RP	particle	<i>up, off</i>			

S1: *Iceland volcano remains hazardous after eruption near Reykjavik.*

S2: *Concerns were raised about the proximity of the volcano to the country's main airport, Keflavik International Airport, which is just a 25 min car ride from the peninsula.*

To annotate the POS tags, you may use the following format:

*John/NNP and/CC Mary/NNP bought/VBD a/DT refrigerator/NN with/IN three/CD doors/NNS*

S1: (2 points)

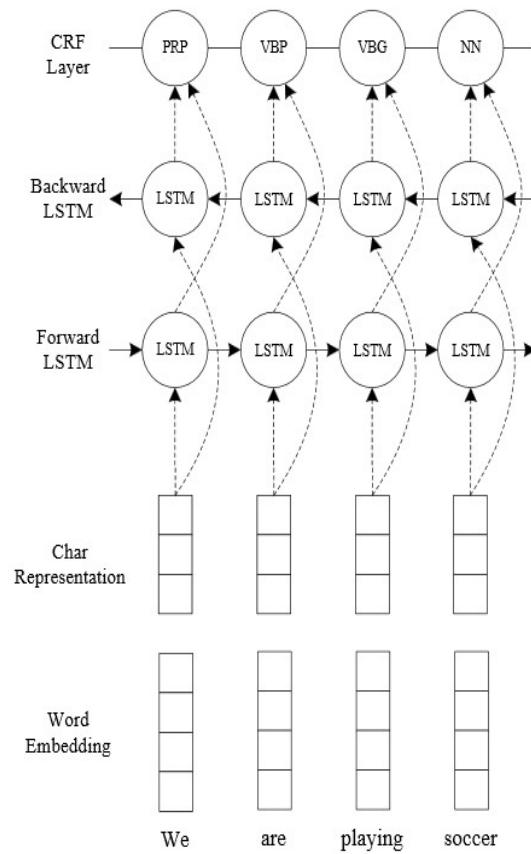
Problem 2 • a) - Solution:  
*S1: Iceland/NNP volcano/NN remains/VBZ hazardous/JJ after/IN  
 eruption/NN ~~at~~ near/IN Reykjavik/NNP.*

S2: (3 points)

S2: concerns/NNS were/VBD raised/VBN about/IN the/DT proximity/NN of/IN the/DT volcano/NN to/IN the/DT ~~country~~/NN 's/pos main/JJ airport/NN, Keflavik/VNP International/VNP Airport/VNP, which/WDT is/VBZ just/RB a/DT 25/CD min/NN car/NN ride/NN from/ZN the/DT peninsula/NN.

b) Neural POS Tagging (10 points)

The architecture of a Neural Conditional Random Field (CRF) used for Part-of-Speech (POS) Tagging is represented below (where word embeddings and character embeddings are concatenated):



Explain how the POS tags are learned to be assigned. **(5 points)** What parameters does this neural architecture have and explain how are they used and learned? **(5 points)**

Question 2 b): Solution:

First, words are transformed into hot vectors by word embeddings and character embeddings.

Then, with bi-directional LSTM layers, we did not compute the probability for each tag at each step. Instead, we used log-linear functions to get a global probability for the whole sequence.

Possible tag sequences:  $\hat{T} = \arg \max_T P(T|w)$

$$P(Y|X) = \frac{\exp\left(\sum_{k=1}^K w_k f_k(X, Y)\right)}{\sum_{Y'} \exp\left(\sum_{k=1}^K w_k f_k(X, Y')\right)}$$

current output token  $y_i$   
 previous output token  $y_{i-1}$   
 input string  $X$   
 current position  $i$

Linear Chain CRF:

$$f_k(X, Y) = \sum_{i=1}^n f_k(y_{i-1}, y_i, X, i)$$

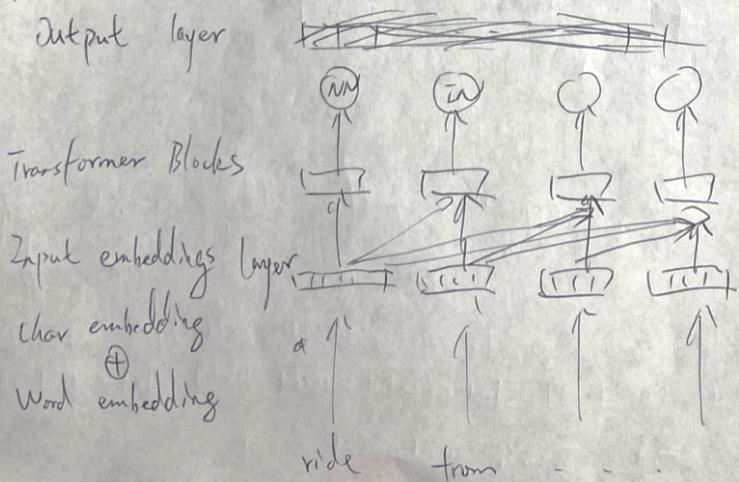
Then, we have:

$$\hat{T} = \arg \max_{T} \sum_{i=1}^n \sum_{k=1}^K w_k f_k(y_{i-1}, y_i, X, i)$$

- c) Propose a transformer-based architecture for performing POS tagging. Draw the architecture, give the equations and discuss how the self-attention is computed **(10 points)**. How do you implement and represent the input to your neural architecture?

**(5 points).**

Problem 2 c) Solution:



Equations and input solution:

$$Q = W^Q X$$

$$K = W^K X$$

$$V = W^V X$$

$$\text{SelfAttention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{dk}}\right)V$$

[connect to Problem 2 c]

Use char embedding and word embedding to firstly transform word into ~~the~~ vectors embedding vectors for input embedding layer.

### Problem 3 (Affect) (20 points)

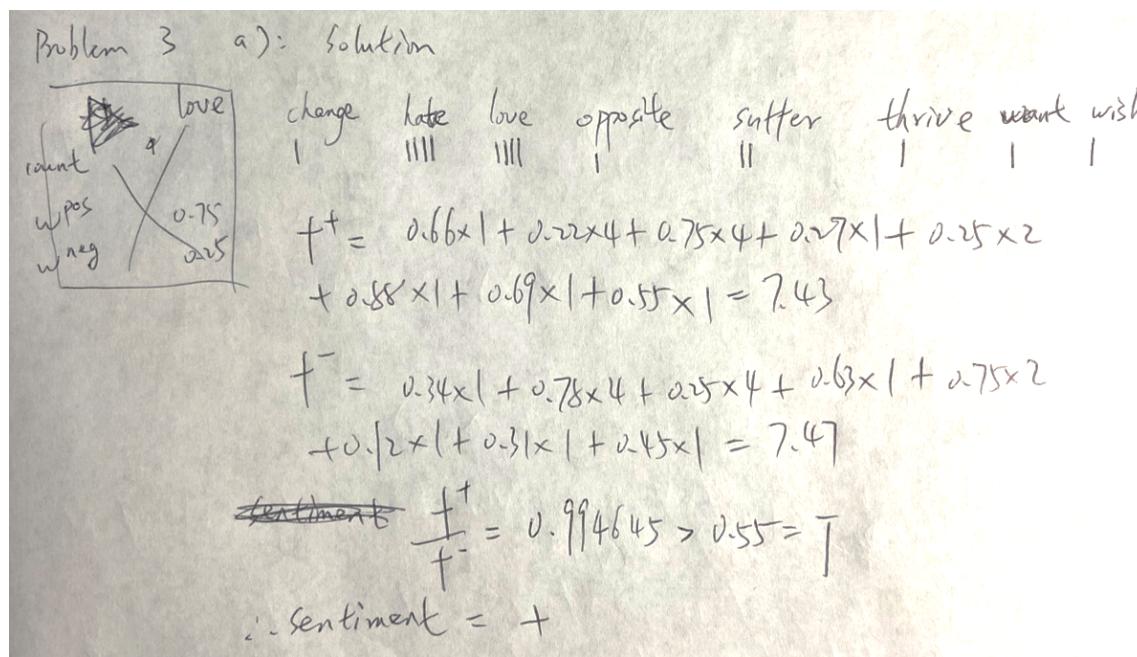
The SentProp Algorithm has computed the following Positive and Negative scores for the polarity of words:

**LEXICON:** { (change:  $w^{pos}=0.66$ ,  $w^{neg}=0.34$ ); (hate:  $w^{pos}=0.22$ ,  $w^{neg}=0.78$ ), (love:  $w^{pos}=0.75$ ,  $w^{neg}=0.25$ ); (opposite:  $w^{pos}=0.27$ ,  $w^{neg}=0.63$ ), (suffer:  $w^{pos}=0.25$ ,  $w^{neg}=0.75$ ); (thrive:  $w^{pos}=0.88$ ,  $w^{neg}=0.12$ ), (want:  $w^{pos}=0.69$ ,  $w^{neg}=0.31$ ), (wish:  $w^{pos}=0.55$ ,  $w^{neg}=0.45$ ) }.

- a) Given this lexicon, determine the sentiment of the following text and explain providing details how you obtained the resulting sentiment, considering that the threshold you will consider in your computation is  $T=0.55$  [10 points]

**TEXT:** Love and hate are similar in being directed toward another person because of who he or she is. Despite this similarity, the two seem like polar opposites. Very often when we love someone, we want them to thrive. When we hate someone, we are more likely to wish they would suffer — or at least change who they are. If she loves you, it doesn't mean she is not capable of hate. If you hate him now and suffer because of it, it does not mean you never loved him.

Solution:



Positive

- b) Whenever any word from the lexicon is within the scope of a negation (e.g. “never loved”), swap the values of the positive and negative weights and recompute the sentiment of the text. Show the details of the computation. [10 points]

Solution:

b) Solution								
	change	hate	love	opposite	suffer	thrive	want	wish
remain score	1	-3	3	1	>	1	1	1
Swap value	0	1	1	0	0	0	0	0
$f^+$	$= 0.64 \times 1 + 0.22 \times -3 + 0.78 \times 1 + 1.75 \times 3 + 0.25 \times 1 + 0.7 \times 1 + 0.25 \times 2 + 0.88 \times 1 + 0.69 \times 1 + 0.5 \times 1 = 7.49$							
$f^-$	$= 0.44 \times 1 + 0.78 \times 3 + 0.22 \times 1 + 0.25 \times 3 + 0.75 \times 1 + 0.63 \times 1 + 0.75 \times 2 + 0.12 \times 1 + 0.31 \times 1 + 0.45 \times 1 = 7.41$							
$f^+ / f^-$	$= 1.010796 > 0.55 = T$							
							$\therefore$ sentiment = +	

Positive