

NLP Midterm

Problem 1 a) : Solution:

MLE n-grams (without smoothing):

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

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

Besides, 2500 sentences \Rightarrow 2500 <s> and 2500 </s>

	<s>	Tom	talks	a	lot	sometimes	</s>
<s>	0.0 x 2500	0.65 x 2500	0.0 x 2500	0.15 x 2500	0.0 x 2500	0.2 x 2500	0.0 x 2500
Tom	0.0 x 1500	0.0 x 1500	0.7 x 1500	0.0 x 1500	0.0 x 1500	0.3 x 1500	0.0 x 1500
talks	0.0 x 2000	0.0 x 2000	0.0 x 2000	0.65 x 2000	0.0 x 2000	0.15 x 2000	0.2 x 2000
a	0.0 x 1000	0.0 x 1000	0.0 x 1000	0.0 x 1000	0.9 x 1000	0.1 x 1000	0.0 x 1000
lot	0.0 x 400	0.2 x 400	0.1 x 400	0.0 x 400	0.0 x 400	0.6 x 400	0.3 x 400
sometimes		0.2					
</s>							

	<s>	Tom	talks	a	lot	sometimes	</s>
<s>	0.0 x 2500	0.65 x 2500	0.0 x 2500	0.15 x 2500	0.0 x 2500	0.2 x 2500	0.0 x 2500
Tom	0.0 x 1500	0.0 x 1500	0.7 x 1500	0.0 x 1500	0.0 x 1500	0.3 x 1500	0.0 x 1500
talks	0.0 x 2000	0.0 x 2000	0.0 x 2000	0.65 x 2000	0.0 x 2000	0.15 x 2000	0.2 x 2000
a	0.0 x 1000	0.0 x 1000	0.0 x 1000	0.0 x 1000	0.9 x 1000	0.1 x 1000	0.0 x 1000
lot	0.0 x 400	0.2 x 400	0.1 x 400	0.0 x 400	0.0 x 400	0.6 x 400	0.3 x 400
sometimes	0.0 x 550	0.2 x 550	0.5 x 550	0.2 x 550	0.0 x 550	0.0 x 550	0.3 x 550
</s>	1 x 2500	0.0 x 1500	0.0 x 1500	0.0 x 2500	0.0 x 2500	0.0 x 2500	0.0 x 2500

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.

$$P_2 = P(\text{sometimes} | \text{<s>}) P(\text{Tom} | \text{sometimes}) P(\text{talks} | \text{Tom}) P(\text{a} | \text{talks}) P(\text{lot} | \text{a}) P(\text{</s>} | \text{lot})$$

$$= 0.2 \times 0.2 \times 0.7 \times 0.65 \times 0.9 \times 0.1 = 0.001638$$

Problem 1 c): Solution:

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

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

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

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

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

$N=16$

~~merge replicated items~~

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

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, $|V| = 25$

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

$\therefore P_L(w_i) = \frac{P_m(w_i)N+1}{N+V}$ ~~$P_m(w_i) = \frac{i}{N}$~~

$P_L(\text{Tom} | \langle s \rangle) = \frac{0.7 \times 2500 + 1}{2500 + 25} = 0.693465$ $P_L(\text{talks} | \text{Tom}) = \frac{0.75 \times 1500 + 1}{1500 + 25} = 0.738361$

$P_L(\text{home} | \text{talks}) = \frac{0.1 \times 1500 + 1}{1500 + 25} = 0.099902$ $P_L(a | \text{home}) = \frac{0.1 \times 1500 + 1}{1500 + 25} = 0.099902$

$P_L(\text{lot} | a) = \frac{0.3 \times 1000 + 1}{1000 + 25} = 0.293658$ $P_L(\langle s \rangle | \text{lot}) = \frac{0.05 \times 400 + 1}{400 + 25} = 0.049412$

$P_L(\langle s \rangle | \langle s \rangle) = \frac{1.0 \times 2500 + 1}{2500 + 25} = 0.990495$ $P_L(\text{home} | \langle s \rangle) = \frac{0.1 \times 2500 + 1}{2500 + 25} = 0.099406$

$P_L(\text{Tom} | \text{home}) = \frac{0.5 \times 1500 + 1}{1500 + 25} = 0.499234$ $P_L(\text{talks} | \text{Tom}) = \frac{0.75 \times 1500 + 1}{1500 + 25} = 0.738361$

$P_L(\langle s \rangle | \text{talks}) = \frac{0.05 \times 2000 + 1}{2000 + 25} = 0.049876$ $P_L(\langle s \rangle | \langle s \rangle) = 0.990495$

$P_L(\text{home} | \langle s \rangle) = 0.099406$ $P_L(\text{talks} | \text{home}) = \frac{0 \times 1500 + 1}{1500 + 25} = 0.000066$

~~$P_L(\langle s \rangle | \text{talks}) = 0.049876$~~ $P_L(\langle s \rangle | \text{talks}) = 0.049876$

$\therefore PP_{\text{Laplace}}(w) = \left[0.693465 \times 0.738361 \times 0.099902 \times 0.099902 \times 0.293658 \times 0.049412 \times 0.990495 \times 0.099406 \times 0.499234 \times 0.738361 \times 0.049876 \times 0.990495 \times 0.099406 \times 0.000066 \times 0.049876 \right]^{\frac{1}{16}} = 6.843318$

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

[connect to Problem 1 c)]

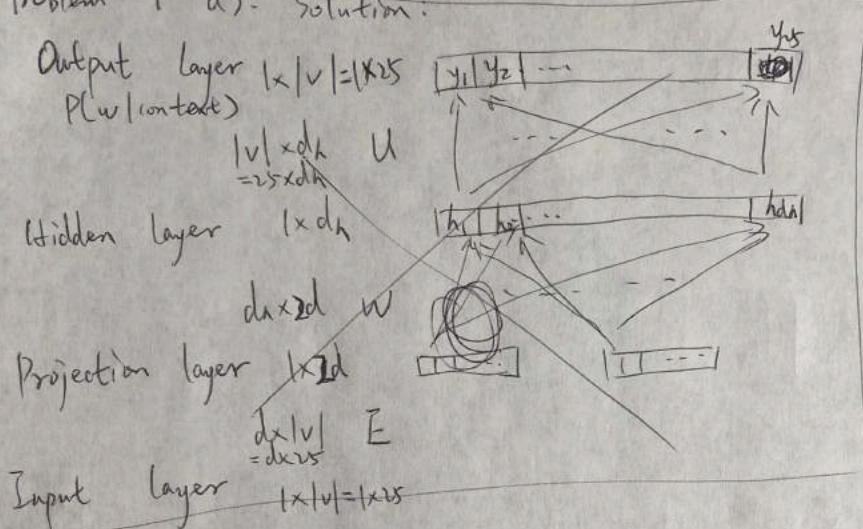
$$PP_{\text{no-smoothing}}(w) = \left[\frac{1}{0} \right]^{\frac{1}{10}} = \infty$$

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

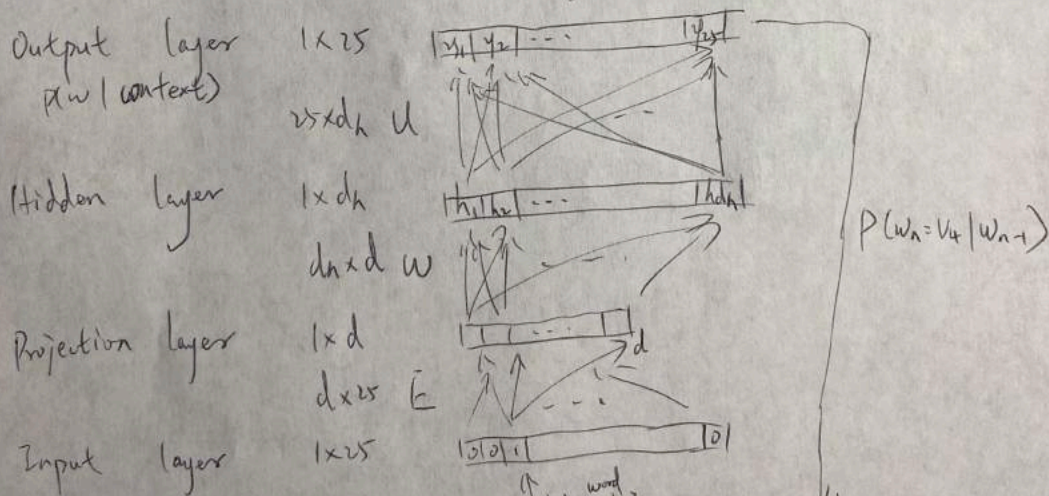
$$PP_{\text{no-smoothing}}(w) = \infty$$

That is why we apply Laplace smoothing for calculating perplexity

Problem 1 d): Solution:



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



<=> When at home Tom speaks

(a)
 w_{n-1}

(b)
 w_n

$$e_i = \bar{e} x_i$$

$$h = \sigma(Wx + b)$$

$$z = Uh$$

$$y = \text{softmax}(z)$$

[connect to Problem 1 d)]

Equations

Problem 1 e: Solution:

Output layer 1×7

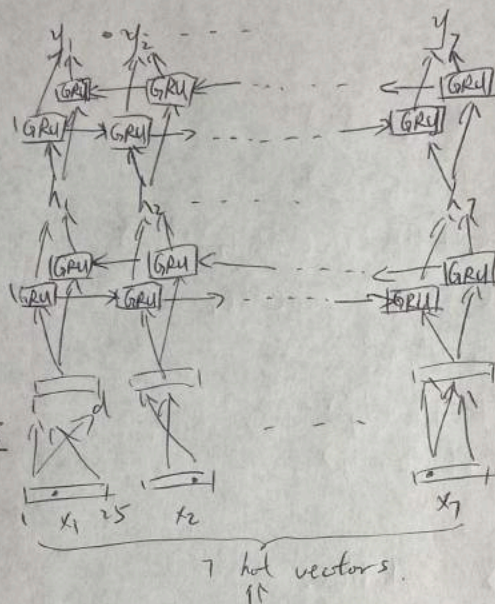
Bi-GRU layer

Hidden layer (x)

Bi-GRU layer

Projection layer 1×7
 $d \times 25 \in$

Input layer $1 \times 25 \times 7$



Lookup	Table					
sometimes	at	home	Tom	speaks	a	lot
w_1	w_2	w_3	w_4	w_5	w_6	w_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.

[connects to Question 1 e]

$$x_i = p_i \times \text{Wembedding}$$

$$r_t = \sigma(u_r x_t + u_r h_{t-1} + b_r)$$

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

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

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

$$h_t^b = w_h^+ h_m^+ + w_h^b h_o^b$$

~~$$h_t^b = \sigma(w_h^+ h_m^+ + w_h^b h_o^b)$$~~

$$v_{t_2} = \sigma(w_{v_2} h_{t_2} + u_{v_2} h_{t_2} + b_{v_2})$$

$$z_{t_2} = \sigma(w_{z_2} h_{t_2} + u_{z_2} h_{t_2} + b_{z_2})$$

$$\hat{h}_{t_2} = \tanh(w_{h_2} h_{t_2} + u_{h_2} (v_{t_2} \odot h_{t_2}) + b_{h_2})$$

$$h_{t_2} = (1 - z_{t_2}) \odot h_{t_2} + z_{t_2} \odot \hat{h}_{t_2}$$

$$e_t = w_e^+ h_m^+ + w_e^b h_o^b$$

$$P(y=j|e) = \frac{\exp(w_j^T e + b_j)}{\sum_{j=1}^K \exp(w_j^T e + b_j)}$$

Problem 2 a) = Solution:

S1: Iceland/NNP volcano/NN remains/VBZ hazardous/JJ after/IN eruption/NN ~~at~~ near/IN Reykjavik/NNP.

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

Question 2 b): Solution:

First, words are transformed into hot vectors by word embeddings and characters 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: $\hat{T} = \underset{T}{\operatorname{argmax}} P(T/w)$

$$P(Y/x) = \frac{\exp(\sum_{k=1}^K w_k f_k(X, Y))}{\sum_{Y \in \mathcal{Y}} \exp(\sum_{k=1}^K w_k f_k(X, Y))}$$

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

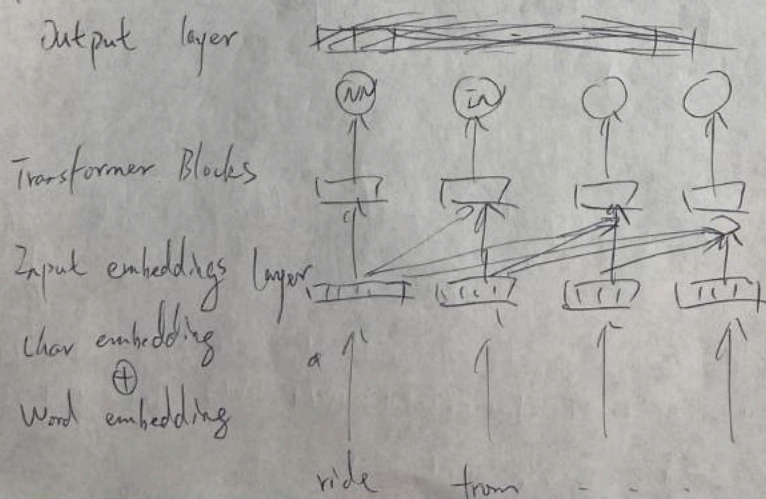
Linear chain CRF:

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

Then, we have:

$$\hat{T} = \underset{Y \in \mathcal{Y}}{\operatorname{argmax}} \sum_{i=1}^n \sum_{k=1}^K w_k f_k(y_{i-1}, y_i, x, i)$$

Problem 2 c): Solution:



[connect to Problem 2 c]

$$Q = W^Q x$$

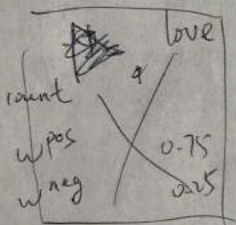
$$K = W^K x$$

$$V = W^V x$$

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

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

Problem 3 a): Solution



change hate love opposite suffer thrive want wish
| |||| | || | || | | |

$$f^+ = 0.66 \times 1 + 0.22 \times 4 + 0.75 \times 4 + 0.27 \times 1 + 0.25 \times 2 + 0.88 \times 1 + 0.69 \times 1 + 0.55 \times 1 = 7.43$$

$$f^- = 0.34 \times 1 + 0.78 \times 4 + 0.25 \times 4 + 0.63 \times 1 + 0.75 \times 2 + 0.12 \times 1 + 0.31 \times 1 + 0.45 \times 1 = 7.47$$

$$\frac{f^+}{f^-} = 0.994645 > 0.55 = T$$

$\therefore \text{sentiment} = +$

b): Solution

change hate love opposite suffer thrive want wish
remain same 1 3 3 1 2 1 1 1
Swap value 0 1 1 0 0 0 0 0

$$f^+ = 0.66 \times 1 + 0.22 \times 3 + 0.78 \times 1 + 0.75 \times 3 + 0.25 \times 1 + 0.27 \times 1 + 0.25 \times 2 + 0.88 \times 1 + 0.69 \times 1 + 0.55 \times 1 = 7.19$$

$$f^- = 0.34 \times 1 + 0.78 \times 3 + 0.22 \times 1 + 0.25 \times 3 + 0.75 \times 1 + 0.63 \times 1 + 0.75 \times 2 + 0.12 \times 1 + 0.31 \times 1 + 0.45 \times 1 = 7.41$$

$$\frac{f^+}{f^-} = 1.010796 > 0.55 = T \quad \therefore \text{sentiment} = +$$