

CS 7650 : Midterm

1) WORD EMBEDDINGS

$$(1) \quad \bar{v}_m = \frac{1}{2h} \sum_{n=1}^h (v_{w_{m+n}} + v_{w_{m-n}})$$

$$\bar{v}_m = \frac{1}{2(3)} \sum_{n=1}^3 (v_{w_{m+n}} + v_{w_{m-n}})$$

$$\bar{v}_3 = \frac{1}{6} \cdot \sum_{n=1}^3 (v_{w_{3+n}} + v_{w_{3-n}})$$

$$= \frac{1}{6} \cdot [v_{w_4} + v_{w_2} + v_{w_5} + v_{w_1} + v_{w_6} + v_{w_0}]$$

$$= \frac{1}{6} [[5, -1] + [-2, 0] + [2, -3] + [3, 2] + [2, 1] + [2, 1]]$$

$$= \frac{1}{6} \cdot [12, 0]$$

$$\boxed{\bar{v}_3 = [2, 0]}$$

(2) word	embedding	unnormalized score	normalized score
dog	[1, 2]	2	0.010
horse	[3, 4]	6	0.165
motorcycle	[0, -1]	0	0.003
leopard	[3, 0]	6	0.165
wolf	[4, 0]	8	0.658

I have

Here, ~~is~~ used these formulae to find the scores:

$$\begin{aligned} - \text{unnormalized score} &= u_i \cdot v_j \\ - \text{normalized score} &= \frac{2^{z_i}}{\sum_j 2^{z_j}} \end{aligned}$$

(3) The CBOW algorithm would predict "wolf" as the missing word, since it has the maximum normalized score.

SHOW WORK

$$[2, 0] \begin{bmatrix} 1 \\ 2 \end{bmatrix} = 2$$

$$[2, 0] \begin{bmatrix} 3 \\ 4 \end{bmatrix} = 6$$

$$[2, 0] \begin{bmatrix} 0 \\ -1 \end{bmatrix} = 0$$

$$[2, 0] \begin{bmatrix} 3 \\ 0 \end{bmatrix} = 6$$

$$[2, 0] \begin{bmatrix} 4 \\ 0 \end{bmatrix} = 8$$

$$\frac{2^2}{389} = 0.010$$

$$\frac{2^6}{389} = 0.165$$

$$\frac{1}{389} = 0.003$$

$$\frac{2^6}{389} = 0.165$$

$$\frac{2^8}{389} = 0.658$$

a.) LSTMs

$$(1) \quad \begin{aligned} f_{m+1} &= \sigma(\theta^{h \rightarrow f} h_m + \theta^{x \rightarrow f} x_{m+1} + b_f) \\ i_{m+1} &= \sigma(\theta^{h \rightarrow i} h_m + \theta^{x \rightarrow i} x_{m+1} + b_i) \\ \tilde{c}_{m+1} &= \tanh(\theta^{h \rightarrow c} h_m + \theta^{x \rightarrow c} x_{m+1}) \end{aligned}$$

LSTM equations

$$c_{m+1} = f_{m+1} \odot c_m + i_{m+1} \odot \tilde{c}_{m+1}$$

$$o_{m+1} = \sigma(\theta^{h \rightarrow o} h_m + \theta^{x \rightarrow o} x_{m+1} + b_o)$$

$$h_{m+1} = o_{m+1} \odot \tanh(c_{m+1})$$

Elman RNN

$$h_t = \tanh(w x_t + V h_{t-1} + b_t)$$

$$y_t = \tanh(u h_t + b_y)$$

By setting $f_{m+1} = 0$, $i_{m+1} = 1$ and $o_{m+1} = 1$, the hidden state at the next time step will only depend on the previous hidden state and the ~~next~~ current input. This effectively turns the LSTM into an Elman RNN.

h_{m+1} only depends on h_m and x_{m+1} , and it doesn't depend on c_m .

(2) \tilde{c}_{m+1} (the candidate memory update) in an LSTM) most closely resembles the hidden state of a standard Elman RNN (h_t).

Both multiply the previous hidden state and the current input with parameters, sum them up (\tilde{c}_{m+1} also adds a bias 'bi') and take the $\tanh()$ of this.

(3) In order to create vector gates, the parameters $\theta(h \rightarrow f)$, $\theta(x \rightarrow f)$, b_f , $\theta(h \rightarrow i)$, $\theta(x \rightarrow i)$, b_i , $\theta(h \rightarrow o)$, $\theta(x \rightarrow o)$, b_o would have to be changed. The following table shows what the dimensions of these parameters would have to be.

Parameter	Dimensions
$\theta(h \rightarrow f)$	$(3, 3)$
$\theta(x \rightarrow f)$	$(3, 3)$
b_f	$(3, 1)$
$\theta(h \rightarrow i)$	$(3, 3)$
$\theta(x \rightarrow i)$	$(3, 3)$
b_i	$(3, 1)$
$\theta(h \rightarrow o)$	$(3, 3)$
$\theta(x \rightarrow o)$	$(3, 3)$
b_o	$(3, 1)$

The benefit of having vector gates instead of scalars is that each "unit" in c_n , \tilde{c}_{n+1} and $\tanh(c_{n+1})$ can be assigned a separate "weight". Using scalar gates imposes a constraint that the "weight" applied to each "unit" must be the same.

As a result, having vector gates gives us the ability to forget certain units, while remembering other units in the same vector. This will allow the LSTM to learn better.

- (4) The inclusion of the memory cell c_{m+1} improves the following problems in RNNs:
- (i) Vanishing gradients: gradients won't decay to 0 during backprop as fast.
 - (ii) Exploding gradients: gradients won't explode to ∞ during backprop as fast.
 - (iii) ~~long~~ ^{short} "memory": information can propagate in the hidden state across several time steps, without attenuating a lot.

The vanishing gradients problem is alleviated since a recurrence function isn't repeatedly applied while computing c_{m+1} .

The short memory problem is alleviated since we use forget and input gates to determine how much of ~~the~~ previous cell state should contribute to the next memory cell state.

3) Beam search Decoding.

(1) Prefix string: "<s> I know"
 $p(\text{prefix}) = 1$, $k=2$

↳ 1st decoding step.

only 2 options: ',' and 'the'

$$\begin{aligned} \cdot \quad p(\text{<s> I know,}) &= p(\text{<s> I know}) \cdot p(, | \text{know}) \\ &= 1 \times 0.4 \\ &= 0.4 \end{aligned}$$

$$\begin{aligned} \cdot \quad p(\text{<s> I know the}) &= p(\text{<s> I know}) \cdot p(\text{the} | \text{know}) \\ &= 1 \times 0.6 \\ &= 0.6 \end{aligned}$$

$$\therefore \text{beam} = [\text{<s> I know, } \underset{0.4}{\text{,}}, \text{<s> I know the} \underset{0.6}{\text{the}}]$$

↳ 2nd decoding step.

$$\begin{aligned} \checkmark \cdot \quad p(\text{<s> I know, I}) &= p(\text{<s> I know,}) \cdot p(I | ,) \\ &= 0.4 \times 0.8 = 0.32 \end{aligned}$$

$$\begin{aligned} \cdot \quad p(\text{<s> I know, it}) &= p(\text{<s> I know,}) \cdot p(it | ,) \\ &= 0.4 \times 0.2 = 0.08 \end{aligned}$$

$$\begin{aligned} \checkmark \cdot \quad p(\text{<s> I know the correct}) &= p(\text{<s> I know the}) \cdot p(\text{correct} | \text{the}) \\ &= 0.6 \times 0.5 = 0.3 \end{aligned}$$

$$\begin{aligned} \cdot \quad p(\text{<s> I know the important}) &= p(\text{<s> I know the}) \cdot p(\text{important} | \text{the}) \\ &= 0.6 \times 0.3 = 0.18 \end{aligned}$$

$$P(<S> \text{ I know the exact}) = P(<S> \text{ I know the}) \cdot P(\text{exact} | \text{the})$$

$$= 0.6 \times 0.2 = 0.12$$

$$\therefore \text{beam} = ["<S> \text{ I know, I}" \downarrow 0.32, "<S> \text{ I know the correct}" \downarrow 0.30]$$

\hookrightarrow 3rd decoding step.

$$\cdot P(<S> \text{ I know, I will}) = P(<S> \text{ I know, I}) \cdot P(\text{will} | \text{I})$$

$$\cdot \cancel{P(<S> \text{ I know})} = 0.32 \times 0.4 = 0.128$$

$$\checkmark \cdot P(<S> \text{ I know, I know}) = P(<S> \text{ I know, I}) \cdot P(\text{know} | \text{I})$$

$$= 0.32 \times 0.6 = 0.192$$

$$\checkmark \cdot P(<S> \text{ I know the correct response}) = P(<S> \text{ I know the correct}) \cdot P(\text{response} | \text{correct})$$

$$= 0.3 \times 0.6$$

$$= 0.18$$

$$\cdot P(<S> \text{ I know the correct answer})$$

$$= P(<S> \text{ I know the correct}) \cdot P(\text{answer} | \text{correct})$$

$$= 0.3 \times 0.25$$

$$= 0.075$$

$$\cdot P(<S> \text{ I know the correct problems}) =$$

$$P(<S> \text{ I know the correct}) \cdot P(\text{problem} | \text{correct})$$

$$= 0.3 \times 0.15 = 0.045$$

$$\therefore \text{beam} = ["<S> \text{ I know, I know}" \downarrow 0.192, "<S> \text{ I know the correct response}" \downarrow 0.18]$$

Most likely sequence: $<S> \text{ I know, I know}$

(a) number of output sequences that can be generated = 7

(b) The sequence with the maximum probability in S is " $\langle s \rangle$ I know, I know" with probability 0.192

(c) The sequence with the minimum probability in S is " $\langle s \rangle$ I know the correct answer" with probability 0.075

t=1

$P(\langle s \rangle I \text{ know}) \cdot P(\text{word} | \text{know})$
 \swarrow
 $I \rightarrow 0.4$
 \searrow
 $\text{the} \rightarrow 0.6$

t=2

$P(\langle s \rangle I \text{ know}) \cdot P(\text{word} | I)$
 \swarrow
 $I \rightarrow 0.32$
 \searrow
 $\text{it} \rightarrow 0.08$

$P(\langle s \rangle I \text{ know, it}) \cdot P(\text{word} | \text{it})$
 \swarrow
 $\text{was} \rightarrow 0.08$

t=3

$P(\langle s \rangle I \text{ know, I}) \cdot P(\text{word} | I)$
 \swarrow
 $\text{will} \rightarrow 0.128$
 \searrow
 $\text{know} \rightarrow 0.192$

$P(\langle s \rangle I \text{ know the}) \cdot P(\text{word} | \text{the})$
 \swarrow
 $\text{exact} \rightarrow 0.12$
 \searrow
 $\text{correct} \rightarrow 0.3$
 \swarrow
 $\text{important} \rightarrow 0.18$

$P(\langle s \rangle I \text{ know the exact}) \cdot P(\text{word} | \text{exact})$

$\text{answer} \rightarrow 0.12$

(2)

(a) prefix string: " $\langle s \rangle I \text{ know}$ "
 $P(\langle s \rangle I \text{ know}) = 1$
 $K = 2$

$P(\langle s \rangle I \text{ know the correct}) \cdot P(\text{word} | \text{correct})$

$\text{response} \rightarrow 0.18$
 $\text{answers} \rightarrow 0.075$
 $\text{problems} \rightarrow 0.045$
 $\text{DON'T CONSIDER SINCE TOPK}$
 $\rightarrow \text{response} \rightarrow 0.09$
 $\rightarrow \text{answers} \rightarrow 0.04$

4) Evaluation

Predicted NER tags
Gold NER tags

(1) Barack and Michele Obama attend the WHCD event at the Hilton Hotel in Washington.

0	0	B-PER	I-PER	0	0	B-LOC
B-PER	0	B-PER	I-PER	0	0	0
0	0	0	B-LOC	0	0	B-LOC
0	0	0	B-LOC	I-LOC	0	B-LOC

True Positives : 2 (Michele Obama, Washington)
 False Positives : 2 (WHCD, Hilton)
 False Negatives : 2 (Barack, Hilton Hotel)
 True Negatives : 1

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{2}{4} = 0.5$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{2}{4} = 0.5$$

$$\text{F1-score} = \frac{2 \cdot p \cdot r}{(p + r)} = 0.5$$

(2)

n-gram precision

$$p_n = \frac{\sum_{n\text{-grams} \in c} \min \left(\max_{i=1 \rightarrow K} \frac{\text{count}_{r_i}(n\text{-gram})}{\text{count}_c(n\text{-gram})} \right)}{\sum_{n\text{-grams} \in c} \text{count}_c(n\text{-gram})}$$

$$BP = \begin{cases} 1 & \text{if } \text{len}(c) \geq \min_{i=1 \rightarrow K} \text{len}(r_i) \\ \exp \left(1 - \frac{\text{len}(c)}{\text{len}(c)} \right) & \text{otherwise} \end{cases}$$

$$BLEU = BP \exp \left[\sum_{n=1}^N (w_n \cdot \ln p_n) \right]$$

$$\text{where } \sum_{n=1}^N w_n = 1$$

c_1 : fruits are good for health

c_2 : vegetables are very important for health

r_1 :

eating fruits is good for health

r_2 : fruits and vegetables are essential for good health.

$$w_1, w_2, w_3 = 1/3, w_4 = 0$$

→ ~~P₁~~
- computing the BLEU score for c_1

$$\begin{aligned} \rightarrow p_1 &= \frac{\min(1, 1) + \min(1, 1) + \min(1, 1) + \min(1, 1) + \min(1, 1)}{5} \\ &= 1 \end{aligned}$$

$$\begin{aligned} \rightarrow p_2 &= \frac{\min(0, 1) + \min(0, 1) + \min(1, 1) + \min(1, 1)}{4} \\ &= 1/2 \end{aligned}$$

$$\begin{aligned} \rightarrow p_3 &= \frac{\min(0, 1) + \min(0, 1) + \min(1, 1)}{3} \\ &= 1/3 \end{aligned}$$

$$\rightarrow \text{len}(c) = 5, \text{len}(r) = 6$$

$$\begin{aligned} \text{BP} &= \exp\left(1 - \frac{6}{5}\right) = \exp(-0.2) \\ &= 0.818 \end{aligned}$$

$$\begin{aligned} \rightarrow \text{BLEU} &= 0.818 \cdot \exp\left[\frac{1}{3} \cdot \ln(1) + \frac{1}{3} \cdot \ln(0.5) + \frac{1}{3} \cdot \ln(1/3)\right] \\ &= 0.818 \cdot e^{-0.597} \\ &= 0.818 \times 0.55 \end{aligned}$$

$$\text{BLEU} = \boxed{0.45016}$$

- computing the BLEU score for c_2

$$\begin{aligned} \hookrightarrow p_1 &= \frac{\min(1,1) + \min(1,1) + \min(0,1) + \min(0,1) + \min(1,1) + \min(1,1) + \min(1,1)}{7} \\ &= 5/7 \end{aligned}$$

$$\begin{aligned} \hookrightarrow p_2 &= \frac{\min(1,1) + \min(0,1) + \min(0,1) + \min(0,1) + \min(1,1) + \min(1,1)}{6} \\ &= 1/2 \end{aligned}$$

$$\begin{aligned} \hookrightarrow p_3 &= \frac{\min(0,1) + \min(0,1) + \min(0,1) + \min(0,1) + \min(1,1)}{5} \\ &= 1/5 \end{aligned}$$

$$\begin{aligned} \hookrightarrow \text{len}(c) &= 7, \text{len}(r) = 6 \\ \text{BP} &= 1 \end{aligned}$$

$$\begin{aligned} \hookrightarrow \text{BLEU} &= 1 \cdot \exp \left[\frac{1}{3} \cdot \ln \frac{5}{7} + \frac{1}{3} \cdot \ln \frac{1}{2} + \frac{1}{3} \cdot \ln \frac{1}{5} + 0 \right] \\ &= e^{-0.879} \\ &= \boxed{0.4149} \end{aligned}$$

5) TRANSFORMER SELF ATTENTION

$$(1) \quad \begin{aligned} q_i &= w^Q x_i \\ k_i &= w^K x_i \\ v_i &= w^V x_i \end{aligned}$$

$$\text{score}(x_i, x_j) = \frac{q_i \cdot k_j}{\sqrt{d_k}}$$

$$\alpha_{i,j} = \text{softmax}(\text{score}(x_i, x_j)) \quad \forall j \leq i$$

$$y_i = \sum_{j \leq i} \alpha_{i,j} \cdot v_j$$

Attention is all you need.

Attention	-2.25	-1.917	1.541	6.0417	4.8125
is	-0.625	-0.625	0.625	2.083	1.5625
all	-2.9375	-1.77	0.458	4.6875	4.5
you	-2.5	-0.833	-1.04	1.0417	2.1875
need	-0.9375	0.229	-1.541	-1.9791	-0.5

NORMALIZED SCORES

WORK

$X \rightarrow (N, D)$ matrix of word embeddings

$$Q = XW^Q$$

$$K = XW^K$$

$$V = XW^V$$

$$\text{normalized scores} = \frac{QK^T}{48}$$

(2)

	Attention	is	all	you	need.
Attention	1.92×10^{-4}	2.68×10^{-4}	8.51×10^{-3}	7.66×10^{-1}	2.24×10^{-1}
is	3.4×10^{-2}	3.4×10^{-2}	1.18×10^{-1}	5.1×10^{-1}	3.63×10^{-1}
all	2.64×10^{-1}	8.49×10^{-4}	7.89×10^{-3}	5.41×10^{-1}	4.49×10^{-1}
you	6.5×10^{-3}	3.44×10^{-2}	2.79×10^{-2}	2.24×10^{-1}	7.06×10^{-1}
need	1.50×10^{-1}	4.82×10^{-1}	8.20×10^{-2}	6.29×10^{-2}	2.32×10^{-1}

Attention values = $\text{softmax} \left(\frac{Q \cdot K^T}{48} \right)$

(3)

Attention	11.254	1.759
is	10.729	1.678
all	7.421	3.560
you	2.922	5.488
need	8.082	0.331

attention head output = attention values $\cdot V$

(4)

Attention	-60.952	28.447	-75.822
is	-47.750	53.749	-49.640
all	13.516	56.943	10.976
you	-9.153	24.379	-16.573
need	-27.528	51.277	-46.528

output of self attention layer = (concat output of each head) $\cdot W^O$

6) BYTE PAIR ENCODING

Training corpus: fresh-, french-, ~~strives~~ ^{pries-}

~~step 1:~~

~~corpus:~~
~~1 fresh-~~
~~1 french-~~

~~step 1:~~

~~corpus~~
~~1 f r e n c h~~

• step 1:

- corpus:

1 f r e s h -
 1 f r e n c h -
 1 f r i e s -

merge f, r

- vocabulary:

-, f, r, e, s, h, n,
 c, i

• step 2

- corpus:

1 f r e s h -
 1 f r e n c h -
 1 f r i e s -

merge f, e

- vocabulary:

-, f, r, e, s, h, n, c, i,
 fr

(tie with e, s and h, -)

• step 3

- corpus

1 f r e s h -
 1 f r e n c h -
 1 f r i e s -

merge h, -

- vocab

-, f, r, e, s, h, n, c, i,
 fr, pre

• step 4

- corpus:

1 pre s h-

1 pre n c h-

1 pr i e s -

- vocab:

-, p, r, e, s, h, n, c, i,

pr, pre, h-

NO bigram occurs more than once