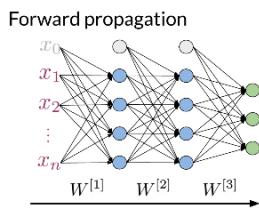


Week 1

Neural Networks with Sentiment Analysis



$$a^{[i]} \text{ Activations ith layer}$$

$$a^{[0]} = X$$

$$z^{[i]} = W^{[i]} a^{[i-1]}$$

$$a^{[i]} = g^{[i]}(z^{[i]})$$

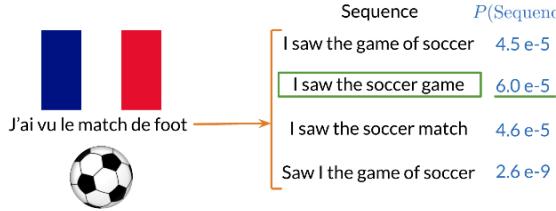
Word	Number
a	1
able	2
about	3
hand	615
...	...
happy	621
...	...
zebra	1000

Tweet: This movie was almost good
 ↓
 [700 680 720 20 55]
 ↓ Padding
 [700 680 720 20 55 | 0 0 0 0 0 0 0]
 To match size of longest tweet

⊕ Train로 구현

RNN for Language Modeling

Traditional LM



- Make use of probabilities to help identify which sentence is most likely to take place

$$P(w_2|w_1) = \frac{\text{count}(w_1, w_2)}{\text{count}(w_1)}$$

$$P(w_3|w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)}$$

$$P(w_1, w_2, w_3) = P(w_1) \times P(w_2|w_1) \times P(w_3|w_2)$$

Bigrams

Trigrams

Bigram 24%.

Large N-grams capture dependencies between distant words and need a lot of space and RAM.

RNN

Nour was supposed to study with me. I called her but she did not

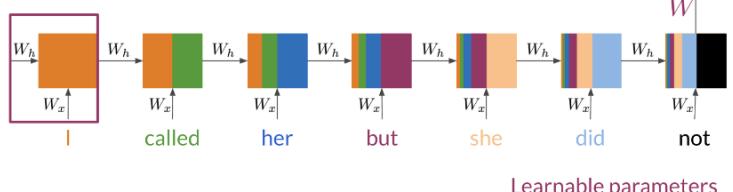
want
respond
choose
want
have
ask
attempt
answer
know

RNNs look at every previous word

Similar probabilities with trigram

Trigram은 2개 단어로 다음 단어 예측

I called her but she did not _____



- 시간이 지남에 따라 단어의 내용 작짐. 그대로 존재!

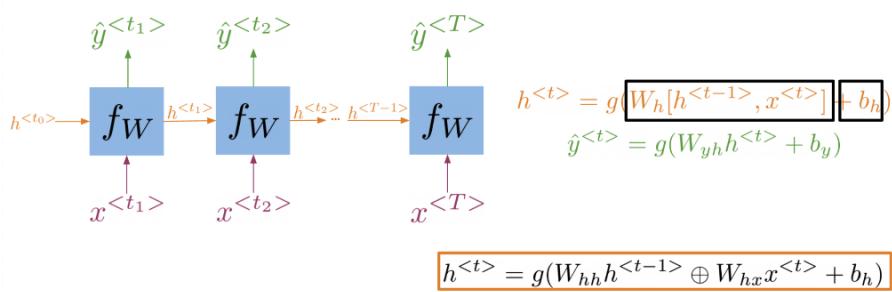
• 많은 계산이 parameter를 공유함

Application of RNNs

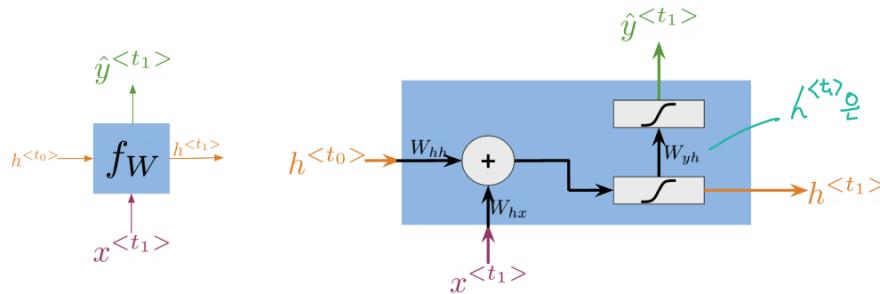
Ex. machine translation, caption generation

- One to One: given some scores of a championship, you can predict the winner.
- One to Many: given an image, you can predict what the caption is going to be.
- Many to One: given a tweet, you can predict the sentiment of that tweet.
- Many to Many: given an English sentence, you can translate it to its German equivalent.

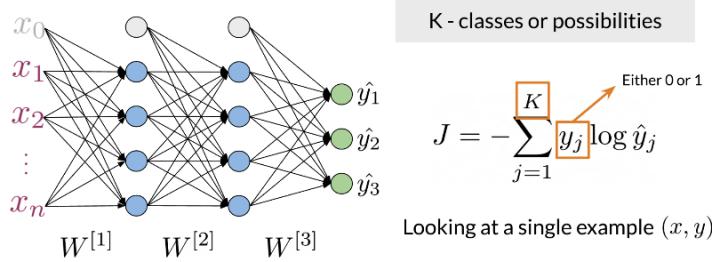
Math in Simple RNNs



$$W_h = (W_{hh}, W_{hx})$$



Cost Function for RNNs



$$J = -\sum_{j=1}^K y_j \log \hat{y}_j$$

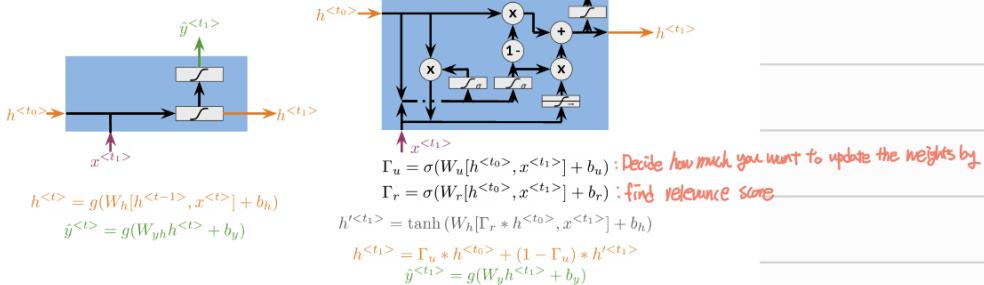
Either 0 or 1

Cross entropy loss

$$\bar{J} = -\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^K y_j^{<t>} \log \hat{y}_{j,t}$$

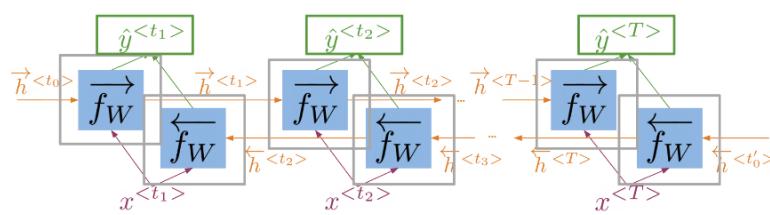
loss over several timesteps

Gated Recurrent Units

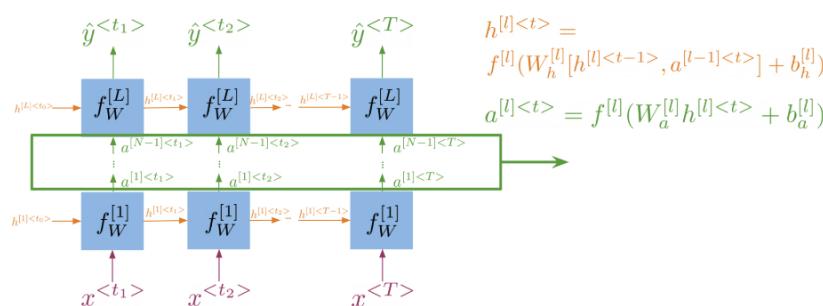


Deep and Bi-directional RNNs

- Bi-directional RNNs: knowing what is next in the sentence could give you more context about the sentence itself



↳ y 만들때 양쪽 방향의 h 사용



Week 3

LSTMs and Named Entity Recognition

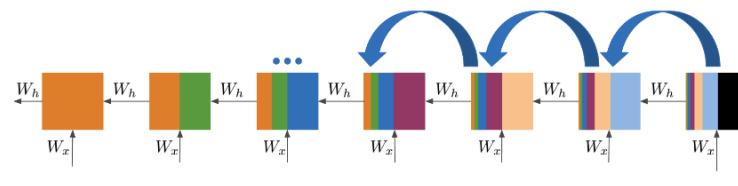
RNNs and Vanishing Gradient

- RNN은 짧은 범위의 dependency를 capture할 수 있고 전통적 LM보다 RAM을 더 사용함

- 단점

① long-term dependency에 취약

② BPTT시에 vanishing/exploding gradient에 취약



$$\text{Sigmoid: } \sigma = \frac{1}{1 + e^{-x}}$$

$$\text{tanh: } \tanh = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

) 이보다 작은 값이 여러 layer를 거치며 곱해지면 결국 가리는 데에 가까워짐

Solutions to Vanishing Gradient Problems

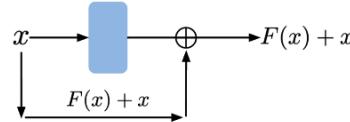
- Identity RNN with ReLU activation

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad -1 \longrightarrow 0$$

- Gradient clipping

$$32 \longrightarrow 25$$

- Skip connections



Introduction To

LSTMs

- cell state + hidden state + 3 gates

Input gate: tells you how much information to input at any time point.

Forget gate: tells you how much information to forget at any time point.

Output gate: tells you how much information to pass over at any time point.

Applications

Next-character prediction

Chatbots



Music composition



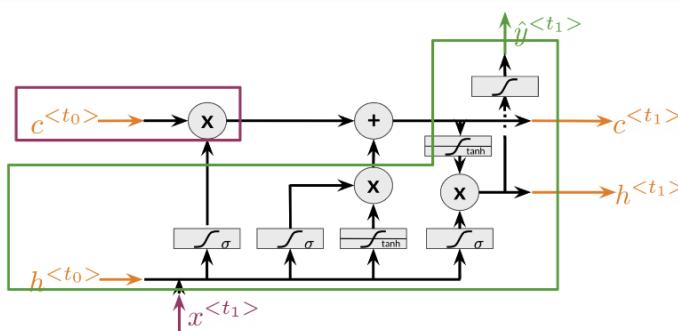
Image captioning



Speech recognition



LSTM Architecture



The forget gate: $f = \sigma(W_f [h_{t-1}; x_t] + b_f)$

The input gate: $i = \sigma(W_i [h_{t-1}; x_t] + b_i)$

The gate gate: $g = \tanh(W_g [h_{t-1}; x_t] + b_g)$

The cell state: $c_t = f \odot c_{t-1} + i \odot g$

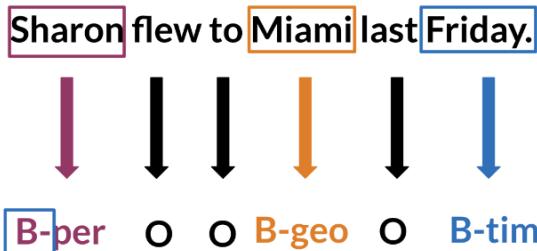
The output gate: $o = \sigma(W_o [h_{t-1}; x_t] + b_o)$

Introduction to

- locates and extracts predefined entities from text.

Named Entity Recognition Ex. places, organizations, names, time, dates

(NER)



Usage : Search efficiency, recommendation engines, customer service, automatic trading

Training NERs :

Data Processing

- Convert words and entity classes into arrays:
- Pad with tokens: Set sequence length to a certain number and use the <PAD> token to fill empty spaces
- Create a data generator:

Sharon flew to Miami last Friday.

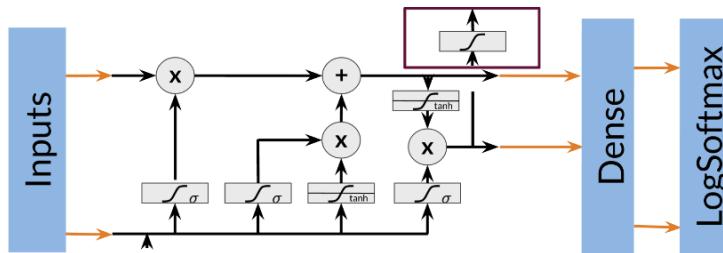
[4282, 853, 187, 5388, 2894, 7]
1 2 3

B-per O O B-geo O B-tim

Training an NER system:

1. Create a tensor for each input and its corresponding number
2. Put them in a batch ==> 64, 128, 256, 512 ...
3. Feed it into an LSTM unit
4. Run the output through a dense layer
5. Predict using a log softmax over K classes

Here is an example of the architecture:



Computing Accuracy

- Pass test set through the model
- Get arg max across the prediction array
- Mask padded tokens
- Compare with the true labels.

Week 4

Siamese Networks

How old are you?

= What is your age?

Where are you from?

≠ Where are you going?

Siamese Networks: learns what makes two inputs the same

Here are a few applications of siamese networks:



Handwritten checks

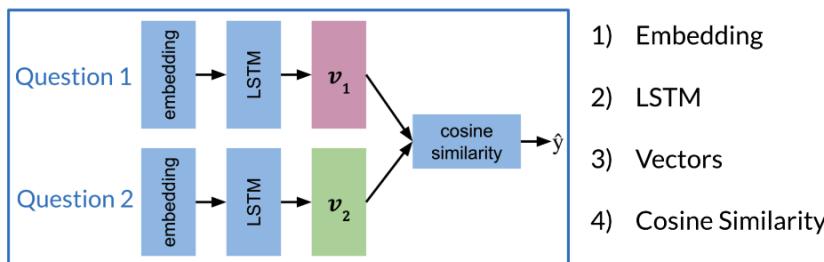
What is your age?
How old are you?



Queries



Question duplicates



• Sub network share identical parameters

→ Need to train one set of weights

Cost function

Triplet loss

How old are you?

Anchor

$$\cos(v_1, v_2) = \frac{v_1 \cdot v_2}{\|v_1\| \|v_2\|}$$

What is your age?

Positive

$$\cos(A, P) \approx 1$$

Where are you from? Negative

$$\cos(A, N) \approx -1$$

$$-\cos(A, P) + \cos(A, N) \leq 0 : \text{We're trying to optimize}$$

$$-2 \leq -\cos(A, P) + \cos(A, N) \leq 2$$

↑ good ↑ bad

Triplets

How old are you?

A

Simplified cost:

$$-\cos(A, P) + \cos(A, N) \leq 0$$

Allows margin of safety

What is your age?

P

$$-\cos(A, P) + \cos(A, N) \leq 0$$

Where are you from?

N

$$-\cos(A, P) + \cos(A, N) \leq 0$$

$$-0.5 + 0.5 + 0.2 \leq 0 \leq 0$$

$$\text{Add a margin: } -\cos(A, P) + \cos(A, N) + \alpha \leq 0$$

$$\text{Full cost: } \max(-\cos(A, P) + \cos(A, N) + \alpha, 0)$$

Do not want negative number as cost

- α : controls how far $\cos(A, P)$ is from $\cos(A, N)$
- **Easy** negative triplet: $\cos(A, N) < \cos(A, P)$
- **Semi-hard** negative triplet: $\cos(A, N) < \cos(A, P) < \cos(A, N) + \alpha$
- **Hard** negative triplet: $\cos(A, P) < \cos(A, N)$

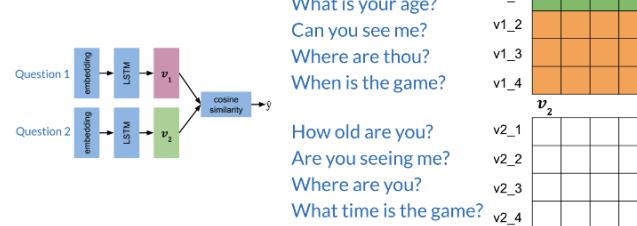
Computing the Cost

Prepare the batches as follows:

What is your age?	How old are you?
Can you see me?	Are you seeing me?
Where are thou?	Where are you?
When is the game?	What time is the game?
$b = 4$	

• 각 row는 같은 의미, 다른 row끼리는 다른 의미

Hard Negative Mining



$v_1 = (1, d_{model})$			
v_{1_1}			
v_{1_2}			
v_{1_3}			
v_{1_4}			

$v_2 = (1, d_{model})$			
v_{2_1}			
v_{2_2}			
v_{2_3}			
v_{2_4}			

• 1 row = 1 sentence's encoding vector

$\cos(v_1, v_2)$				
v_1				
-1	-2	-3	-4	
-1	0.9	-0.8	-0.3	-0.5
-2	-0.8	0.5	-0.1	-0.2
-3	-0.3	-0.1	0.7	-0.8
-4	-0.5	-0.2	-0.8	1.0

$$\text{Cost} = \max(-\cos(A, P) + \cos(A, N) + \alpha, 0)$$

$\cos(v_1, v_2)$				
v_1				
-1	-2	-3	-4	
-1	0.9	-0.8	-0.3	-0.5
-2	-0.8	0.5	-0.1	-0.2
-3	-0.3	-0.1	0.7	-0.8
-4	-0.5	-0.2	-0.8	1.0

mean_neg: speeds up training

closest_neg: helps penalize the cost more

• mean negative of all the other off diagonals in the row

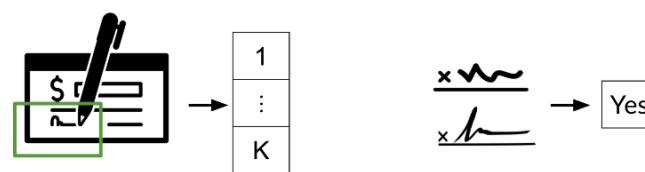
• highest number in the off diagonal

$$\text{Cost1} = \max(-\cos(A, P) + \text{mean}_n eg) + \alpha, 0$$

$$\text{Cost2} = \max(-\cos(A, P) + \text{closest}_n eg + \alpha, 0)$$

The full cost is defined as: Cost 1 + Cost 2.

One shot learning



Classification: classifies input as 1 of K classes

One Shot Learning: measures similarity between 2 classes



Learn a similarity score!

$$\cos(\text{sig1}, \text{sig2}) > \tau \quad \checkmark$$

$$\cos(\text{sig1}, \text{sig2}) \leq \tau \quad \times$$

Testing

Finally when **testing**:

1. Convert two inputs into an array of numbers
2. Feed it into your model
3. Compare v_1, v_2 using cosine similarity
4. Test against a threshold τ