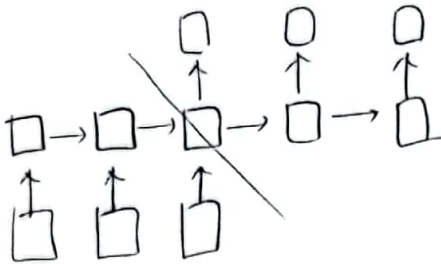
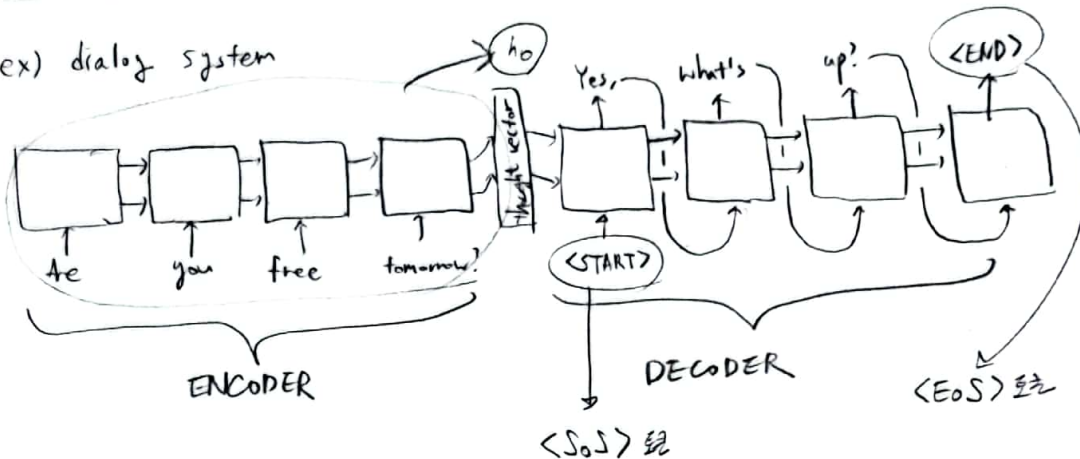


Seq2seq model

RNN 형태 중 many-to-many 모델



ex) dialog system



파라미터를 공유하는 구조 (인코더, 디코더 간)

기본 RNN은 인코더 아키텍처 hidden state에 모든 인코더 정보는 축적된다.
LSTM은 Long-time dependency를 해결했어, 뒤 time step에서는 앞 time step의 정보가 변질되거나 잊혀진다

Attention이 나오기 전에는 입력 순서를 거꾸로 봐도...

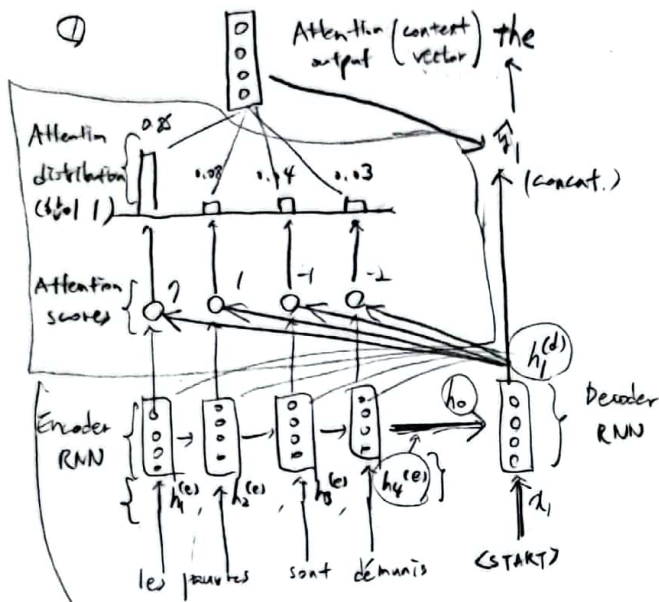
Attention

앞선 예제와, attention을 쓰면

디코더에 인코더 매개변수 time step에 있는 hidden state vector를 제공하는 것이 아니라
각각의 단어는 소문자로 인코딩하는 과정에서 나온 인코딩 hidden state 벡터를 디코더에 제공
 $h_1^{(e)}, h_2^{(e)}, h_3^{(e)}, h_4^{(e)}$

디코더에서는 각 time step에서 단어는 생성한 때, 적절한 인코딩 hidden state 벡터를
선택함으로써 가져가서 예측이 가능

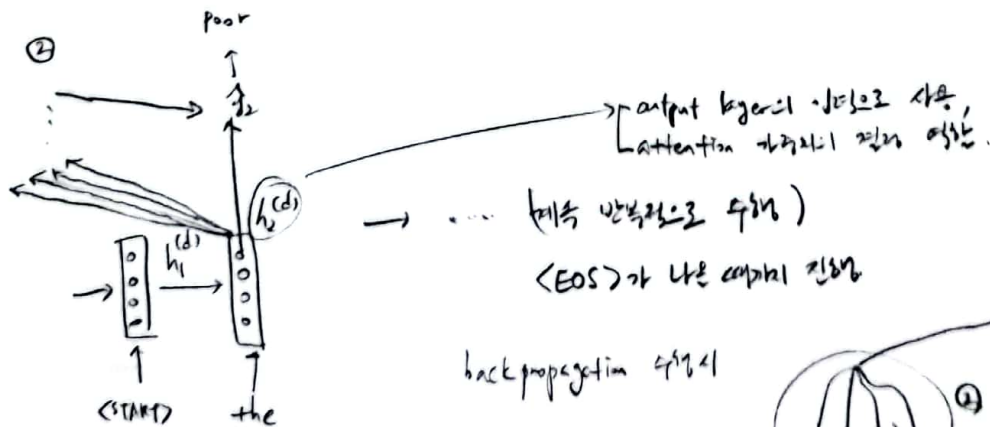
작동 원리



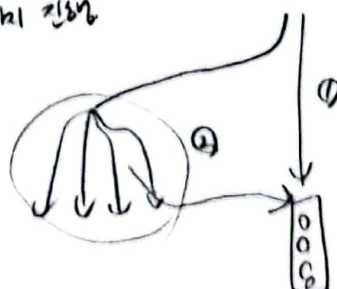
Decoder hidden state 벡터가
Encoder hidden state 각각과 내적 가능
예를 들어 벡터값이 7, 1, -1, -2 라면
(가산)

이것들을 Softmax를 취해서
예를 들어 0.85, 0.08, 0.04, 0.03이 나오면
이것들은 Encoder hidden state 벡터에 부여되는 가중치의 사용
이것을 이용해 가중치들을 계산할 수 있고 이는 뜻이 맞는
해당 Encoding 벡터가 나온다

입력: Decoder hidden state 벡터 사용.
Encoder hidden state 벡터 세트
출력: Encoder hidden state 인코더의 각 time step 벡터 사용



backpropagation 수행시



여기서, decoder의 다음 time step의 입력값으로, ground truth에서의 단어를 사용해야 한다.
 즉, 전 단계에서 예측을 잘못했더라도 그 예측값이 아니라 ground truth에서의 단어를 입력값으로 넣는다
 이런 방식을 Teacher forcing 이라고 한다.

Teacher forcing 이 어떤 때가: 모든 학습을 위해 사용해야 할 때, 성능이 저하된다

(+) 학습이 빠르고 용이하다

(-) Test 할 때의 환경과의 괴리가 있다

이를 보완하기 위해, 하이브리드 Teacher forcing을 적용한다. 어느 정도 학습이 되면
 기준을 넘어서게 되면 된다.

Mechanism

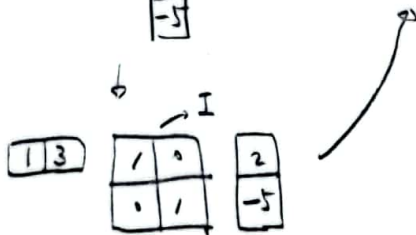
h_t : decoder에서 주어진 hidden state vector

\bar{h}_s : encoder에서 주어진 각 단어의 hidden state vector의 평균

$$\text{Score}(h_t, \bar{h}_s) = \begin{cases} h_t^T \bar{h}_s & \text{Luong - dot} \\ h_t^T W_a \bar{h}_s & \text{Luong - general} \\ v_a^T \tanh(W_a [h_t || \bar{h}_s]) & \text{Luong - concat} \end{cases}$$

↑
가장 보편적

$$\begin{bmatrix} 1 & 3 \end{bmatrix} \begin{bmatrix} 2 \\ -5 \end{bmatrix} = 1 \times 2 + 3 \times (-5)$$



만약

$$\begin{bmatrix} 1 & 3 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 4 \end{bmatrix} \begin{bmatrix} 2 \\ -5 \end{bmatrix}$$

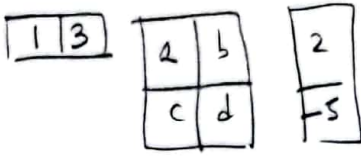
$$\begin{bmatrix} 1 \times 1 & 3 \times 4 \end{bmatrix} \begin{bmatrix} 2 \\ -5 \end{bmatrix} = 1 \times 1 \times 2 + 4 \times 3 \times (-5)$$

각 dimension별로 곱해서 나온 값에 적용하는 계층이

$$\begin{bmatrix} 1 & 3 \end{bmatrix} \begin{bmatrix} 1 & 7 \\ -8 & 4 \end{bmatrix} \begin{bmatrix} 2 \\ -5 \end{bmatrix}$$

$$\begin{bmatrix} 1 \times 1 + 3 \times (-8) & 1 \times 7 + 3 \times 4 \end{bmatrix} \begin{bmatrix} 2 \\ -5 \end{bmatrix} = 1 \times 1 \times 2 + 4 \times 3 \times (-5) + (-8) \times 3 \times 2 + 7 \times 1 \times (-5)$$

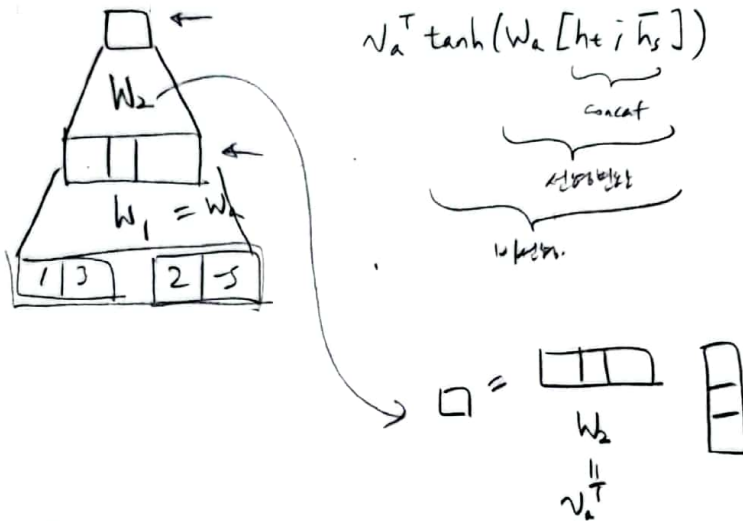
한글



리전 a, b, c, d는 학습 가능한 parameter로 보면

a, b, c, d는 모두 서로 다른 dimension 간의 공역인 값들에 각각 부여되는 것임.

concat은...



원래는 Attention score에 학습 가능 parameter가 있었으나
 여기서 학습 가능 parameter가 포함이 되면,
 backpropagation을 통해 W_a 등은 학습하게 된다.

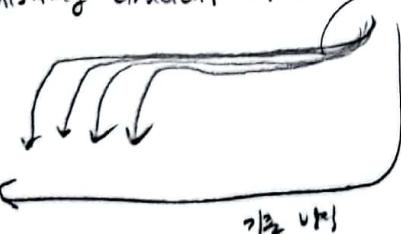
(+) 기계번역 분야에서 성능향상

→ decoder의 매 time step마다 입력 sequence에서 어떤 부분이 정보를 집중시키고 학습하게
 할지 가능하게 했다

bottleneck 현상 해결

Vanishing Gradient 해결

→



Attention 지능
 (어떤 time step과 거리가 있음)

능이로운 해석 가능하게 해줌.