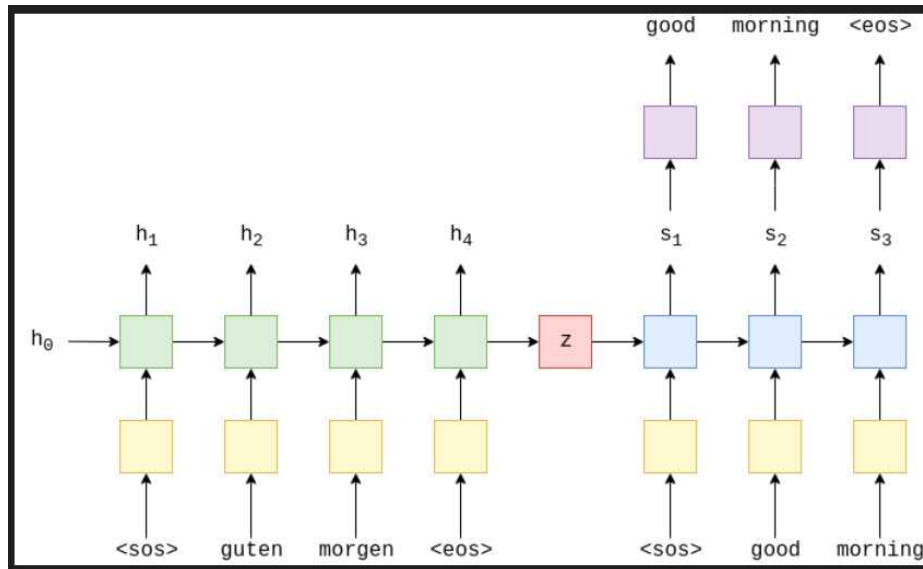


1 - Sequence to Sequence Learning with Neural Networks

DNN은 음성 인식, 사물 인식 등에서 꾸준한 성과를 내어왔다. 하지만 input size가 fixed 된다는 한계점이 존재하기 때문에 sequential problem을 제대로 해결할 수 없다는 한계점이 존재했다. 이 논문에서는 2개의 LSTM을 각각 Encoder, Decoder로 사용해 sequential problem을 해결하고자 했다. 이를 통해 많은 성능 향상을 이루어냈으며, 특히나 long sentence에서 더 큰 상승 폭을 보였다. 이에 더해 단어를 역순으로 배치하는 방식으로 성능을 향상시켰다.

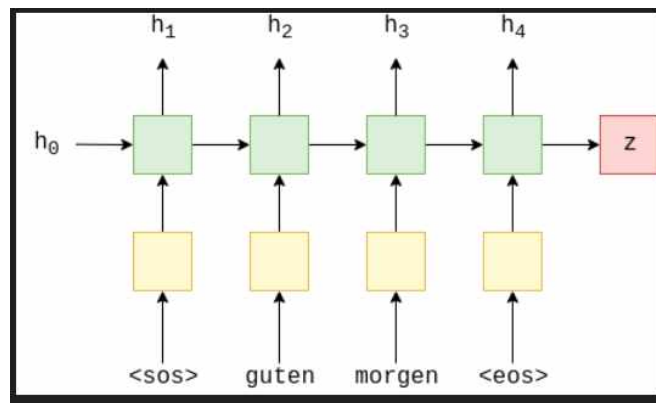


RNN은 사실 sequential problem에 매우 적절한 model이다. 하지만 input size와 output size가 다른 경우에 대해서는 좋은 성능을 보일 수 없었다. 본 논문에서 제시하는 model은 Encoder LSTM에서 하나의 context vector를 생성한 뒤 Decoder LSTM에서 context vector를 이용해 output sentence를 생성하는 방식으로 RNN의 한계점을 극복하고자 했다. input과 output sentence 간의 mapping을 하는 것이 아닌, input sentence를 통해 encoder에서 context vector를 생성하고, 이를 활용해 decoder에서 output sentence를 만들어내는 것이다. Encoder LSTM의 output인 context vector는 Encoder의 마지막 layer에서 나온 output이다. 이를 Decoder LSTM의 첫 번째 layer의 input으로 넣게 된다. 여기서 주목할만한 점은 input sentence에서의 word order를 reverse 해 사용했다는 것이다. 또한 EOS token을 각 sentence의 끝에 추가해 variable length sentence를 다루었다.

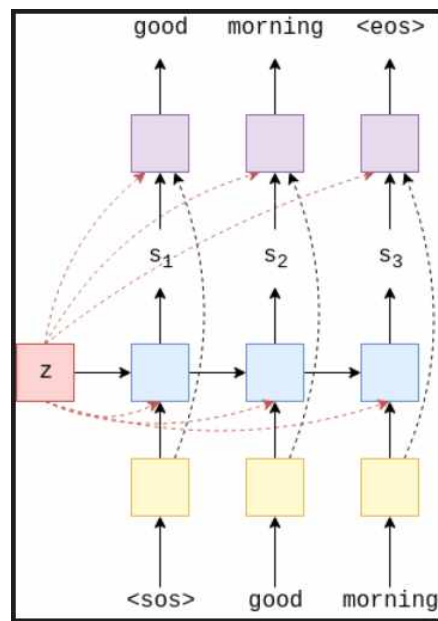
2 - Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation

이 논문은 RNN(GRU)를 이용한 Encoder Decoder 방식을 이용한 Translation Model을 제시함으로써 NMT방식을 제안하고 있다. 다만 이 논문은 NMT라는 용어를 사용하지 않고, 조심스럽게 SMT 이외의 다른 방식이 존재하고, 이러한 방식이 충분히 가능성이 있다는 approach를 보여준다. 이는 SMT방식의 process 중 일부를 대체하는 방식에서 확인할 수 있다.

이 논문에서 제안하는 NMT방식은 seq2seq방식이다. 다만 우리가 흔히 알고 있는 구조와는 다른 방식을 보여준다. 우리가 seq2seq방식이라 하면 Encoding Sequence의 마지막 output EOS와 마지막 Cell State를 Input으로 받아서 Decoder가 처리해주는 방식으로 보여주게 되어 위의 방식보다는 조금 더 많은 정보를 처리하는 것을 볼 수 있다.



Encoder 구조는 흔히 알려져 있는 seq2seq 방식과 크게 차이가 없다.



Decoder 방식은 다르다.

3 - Neural Machine Translation by Jointly Learning to Align and Translate

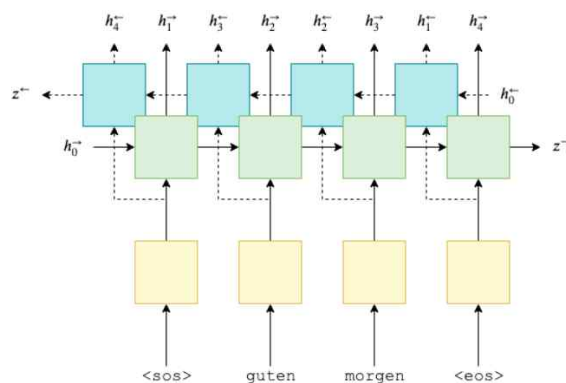
앞서 말한 논문들은 Machine Translation, 그중에서도 NMT를 다루고 있다. 이렇게 소개된 NMT의 구조는 주로 Encoder-Decoder방식을 사용하고 있는데, Encoder의 경우, 입력 문장을 고정 길이 벡터로 변환시키고, decoder는 해당 벡터를 이용해서 번역 결과를 생성해낸다. 해당 모델을 입력으로 제공된 Input sentence와 결과로 생성된 output sentence 사이의 probability를 최대화하는 방식으로 학습된다.

이러한 encoder-decoder 구조는 bottleneck 문제를 발생시킬 수 있다. Bottleneck 문제는 encoder에서 전체 문장을 하나의 고정된 길이의 벡터로 생성할 때 발생한다. 이러한 문제는 문장이 길어질수록 심각하게 나타나고, encoder-decoder 모델의 성능 또한 문장의 길이가 길어질수록 떨어지게 된다. 문장 전체의 정보를 짧은 고정 길이의 벡터로 모두 나타내기 어렵기 때문이다.

이 논문에서는 encoder-decoder 모델에서 새로운 구조를 추가해서 성능을 향상할 수 있는 방법을 제안한다. 해당 모델은 decoder에서 하나의 결과를 만들어낼 때마다, 입력 문장을 순차적으로 탐색해서 현재 생성하려는 부분과 가장 관련있는 영역을 적용시킨다. 최종적으로 encoder에서 생성한 context word 중 관련성이 크다고 판단되는 영역들과, decoder에서 이미 생성한 결과를 기반으로 다음 단어를 결과로 생성해낸다.

이러한 방식의 가장 큰 장점 중 하나는 입력 문장을 고정된 길이의 벡터로 표현하지 않아도 된다는 것이다. Decoder에서 연산을 진행하면서 encoder에서 생성한 context word를 계속해서 참조하기 때문에, 전체 문장의 정보를 하나의 벡터에 담으려고 하지 않아도 되고, 문장의 길이가 길어지더라도 성능을 유지할 수 있다.

align and translate (Attention)

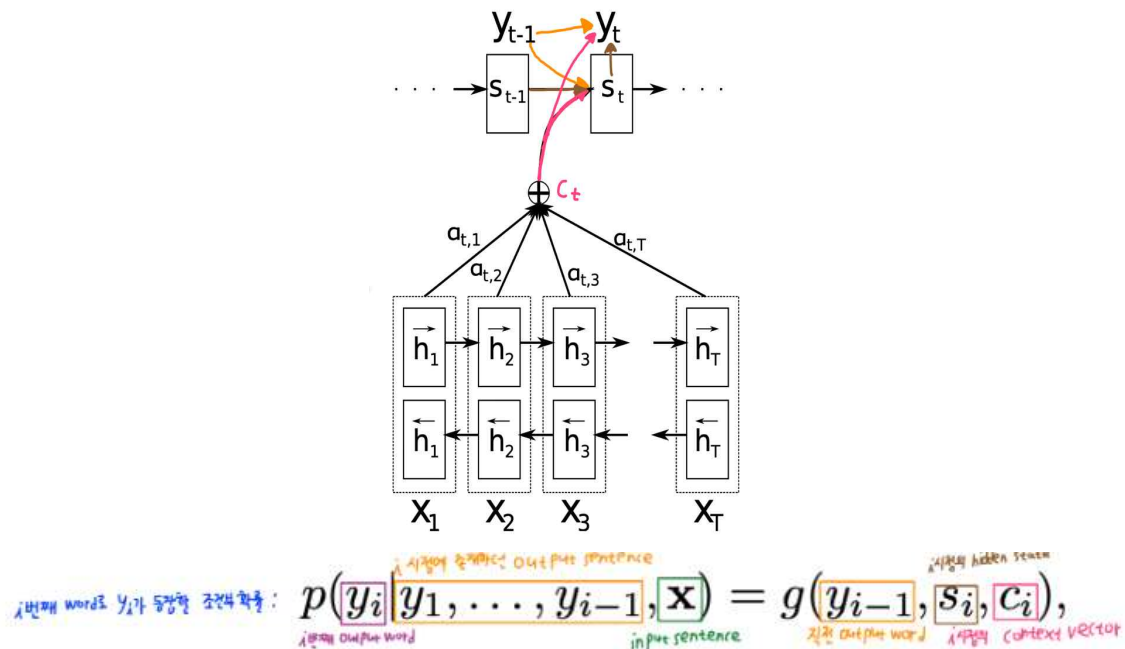


We now have:

$$h_t^{\rightarrow} = \text{EncoderGRU}^{\rightarrow}(e(x_t^{\rightarrow}), h_{t-1}^{\rightarrow})$$
$$h_t^{\leftarrow} = \text{EncoderGRU}^{\leftarrow}(e(x_t^{\leftarrow}), h_{t-1}^{\leftarrow})$$

기존 seq2seq model에서는 encoder는 1개의 context vector를 생성해내기 위해 존재했다. 하지만 본 model은 위에서 언급했듯이 하나의 context vector가 아닌 각 시점에 대한 context vector를 각각 생성한다. 이 때 사용하는 h_j 를 구하는 것이 encoder의 역할이다. h_j 는 위에서 언급했듯이 j 번째 input token에 대해 attention한 vector이다. 이 때 attention한다는 의미는 j 번째 input token에 대해 당연히 가장 높은 가중치를 주고, j 번째에서 멀어질수록 낮은 가중치를 주는 것이다. 이는 $j-1, j-2, \dots$ 의 역방향, $j+1, j+2, \dots$ 의 순방향, 즉 양방향에 대해 모두 적용되어야 한다. 이를 위해 사용한 것이 Bidirectional RNN이다. 순방향에 대한 h_j 를 생성하고, input sentence의 끝에 도달하면 다시 역방향에 대한 h_j 를 생성한다. 그 후 순방향에 대한 $h_j \rightarrow h_j$ 와 역방향에 대한 $h_j \leftarrow h_j$ 를 함께 반영해 최종 h_j 를 만들어낸다.

Decoder



$$\mathbf{s}_i = f(\mathbf{s}_{i-1}, y_{i-1}, \mathbf{c}_i).$$

i 번째 hidden state, $i-1$ 번째 hidden state, i 번째 output word, i 번째 context vector

i 번째 token으로 y_i 가 등장할 조건부 확률에 대한 수식이다. $y_1 \sim y_{i-1}$ (y_i 이전의 output sentence)와 \mathbf{x} (input sentence 전체)에 대해 다음 token으로 y_i 가 생성될 조건부 확률이다. 이는 g 함수에 $y_{i-1}, \mathbf{s}_i, \mathbf{c}_i$ 를 인자로 넣어 생성된 값이다. y_{i-1} 은 직전 시점 $i-1$ 에서 생성한 output token이고, \mathbf{s}_i 는 현재 시점 i 에서의 RNN hidden state, \mathbf{c}_i 는 현재 시점 i 에 생성된 context vector이다. 직관적으로 해석해보면 이전 단어 y_{i-1} 이후에 나올 단어 y_i 를 예측하는 것인데, 이 때 이전 output sentence의 상태 정보를 모두 포함하고 있는 \mathbf{s}_i 를 입력으로 받음으로써 output sentence의 문맥을 반영하고, input sentence에 대한 context vector \mathbf{c}_i 를 통해 input sentence의 문맥을 반영한다. 이전 seq2seq model에서는 context vector가 input sentence 전체에 대한 vector였는데, 이번 attention seq2seq model에서는 특정 시점 i 에 대한 context vector \mathbf{c}_i 가 주어진다. 즉, context vector가 input sentence 전체에 대한 하나의 vector가 아니라 각 시점 i 에 대해 \mathbf{c}_i 가 각각 정의된다는 것이다. 아래에서는 \mathbf{c}_i 에 대해 좀 더 자세하게 살펴본다.

$$\mathbf{c}_i = \sum_{j=1}^T a_{ij} \mathbf{h}_j$$

\mathbf{c}_i 는 a_{ij} 와 \mathbf{h}_j 에 대해 j 부터 T 까지 더한 vector이다. j 부터 T 까지의 의미는 input sentence의 처음부터 끝까지 각 input token에 대해 j 로 순회한다는 의미이다. 즉, j 는 input sentence에서의 token index이다. 반대로 i 는 output sentence에서 현재 시점 i 이다. 정리하자면 j 는 input에 대한 index, i 는 output에 대한 index이다. 우리는 output의 i 시점에서 생성되는 context vector \mathbf{c}_i 에 대한 수식을 살펴보는 것이다. \mathbf{h}_j 는 input sentence 전체의 context를 포함하지만 동시에 특히 j 번째 token과 그 주변에 대해 더 attention(집중)을 한 vector이다. a_{ij} 는 i 번째 output token이 j 번째 input token과 align될 확률값을 뜻한다. a_{ij} 에 대해 더 살펴보자.

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})},$$

e_{ij} 의 softmax
 output의 i 와 input의 j 가 얼마나 서로 match하는가
 $e_{ij} = a(s_{i-1}, h_j)$
 output의 i 번째 hidden state, input의 j 에 focus

e_{ij} 는 i 번째에 들어올 output token과 j 번째 input token이 얼마나 서로 잘 match되는지에 대한 값이다.

output의 문맥을 반영하기 위해 s_{i-1} 를 input으로 받고, j 번째 input token에 대한 attention을 주기 위해 h_j 를 input으로 받는다. 이렇게 완성된 e_{ij} 를 softmax한 α_{ij} 는 i 번째 들어올 output token과 j 번째 input token이 align될 확률값을 뜻한다.

다시 C_i 의 의미로 되돌아와보면, C_i 는 i 시점에서 모든 input sentence token j 에 대해 α_{ij} 와 h_j 를 곱한 vector에 대한 합이다. α_{ij} 는 현재 시점 i 에서 생성될 output token이 j 번째 input token과 align될 확률값이며, h_j 는 input sentence에 대한 context이되, j 번째 input token에 특히 attention한 vector이다. 결론적으로 C_i 는 현재 시점 i 에서 input sentence의 어느 token에 더 attention을 해야 하는지에 대한 context라고 직관적으로 해석 가능하다.

위의 model이 기존 seq2seq model에 비해 가지는 이점은 encoder가 고정된 size context vector 1개에 모든 input sentence token에 대한 attention을 담을 필요가 없다는 것이다. 왜냐하면 decoder model에서 단지 context vector 1개만 사용해 output sentence를 생성해내는 것이 아니라, h_j 와 같은 input sentence token에 대한 attention data를 사용하기 때문이다.