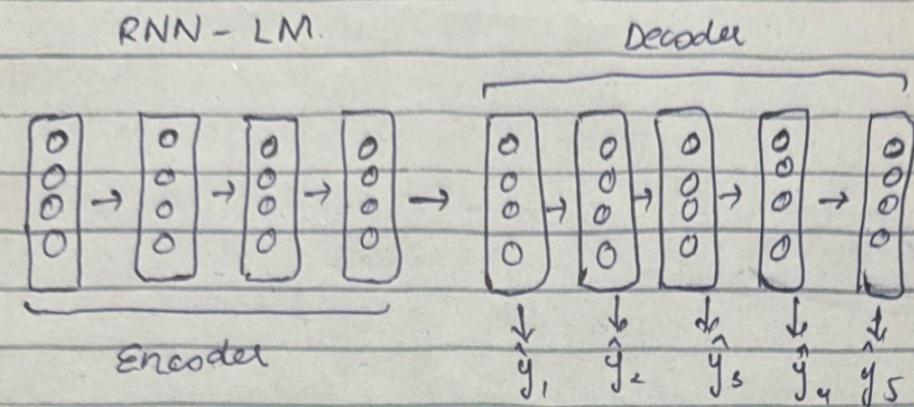


Natural Language Generation (NLG)

↳ generation of new text



For loss : $J_1 = -\log y_1$
 $J_2 = -\log y_2$ & so on

$$\text{Average loss} = \frac{1}{N} \sum J$$

$$= \frac{1}{5} [-\log \hat{y}_1 - \log \hat{y}_2 - \log \hat{y}_3 - \log \hat{y}_4 - \log \hat{y}_5]$$

Decoder Algorithms for generating text

↳ Greedy Decoding

Beam Search

↳ Sampling Based

Greedy Decoding

→ most probable word (argmax)

→ Input Output of one word

is input of next until <END>

↳ Poor output due to lack of backtracking.

(Anggrammatical, Unnatural, Nonsensical,
Incoherent)

search

Beam Decoding

- ↪ Finds highest probability sequences not necessarily the optimal ones.
- ↪ Tracks multiple possible sequence at once.
- ↪ K most probable partial sequences at every step.
 - ↪ $K=1$ ungrammatical, unnatural, incoherent.
 - ↪ K large reduces some of problems but high computation power and more generic output and low BLEU score.

SLG

Beam Search

The

$K = 2$.

The

— Dog.

Bird.

cat

0.26

is
0.5
 0.34×0.5
(0.17) ✓

sat
0.4
 0.34×0.4
(0.13)

barked
0.33
 0.26×0.33
(0.08)

chased
0.37
 0.26×0.37
(0.09)

flew
0.5
 0.4×0.5
(0.2) ✓

saw
0.4
 0.4×0.4
(0.16) ✗

playing
0.5
 0.17×0.5
(0.085)

catching
0.3
 (0.17×0.3)
(0.05)

Away
0.4

Highly
0.95
catching
0.4 ✓

catching
0.34

The cat is playing

The bird flew high.

Sampling Based.

- ↳ Randomness instead of most probable.
- ↳ Types
- ① Pure Sampling.
 - ↳ pick randomly.
 - ↳ e.g. cat 50%, dog 30%, bird 20%.
 - ↳ you pick dog despite not being most probable
- ② Top n sample.
 - ↳ pick any 1 randomly from top n most probable
 - ↳ e.g. cat 50%, dog 30%, bird 20%, fish 10%
 - ↳ Top 2 = cat, dog.
 - ↳ choose dog randomly
 - ↳ keep the value of n small for generic & safe outputs

Decoding

Softmax Temperature (Not decoding algo).

- ↳ Adjust the probabilities on basis of temp
- ↳ T → temp parameter (To control diversity)
- $$p_t(w) = \frac{\exp(Sw/T)}{\sum w'(\exp Sw'/T)}$$
- ↳ Effect of Temp
 - ↳ $T > 1$ (high) (flatten)
uniform probabilities, diverse output
 - ↳ $T < 1$ (low) (sharpened)
concentrated probabilities, less diverse outputs

Example:

$$S_{\text{cat}} = 2, S_{\text{dog}} = 1, S_{\text{bird}} = 0.5$$

$$\text{Temperature } (T) = 0.5$$

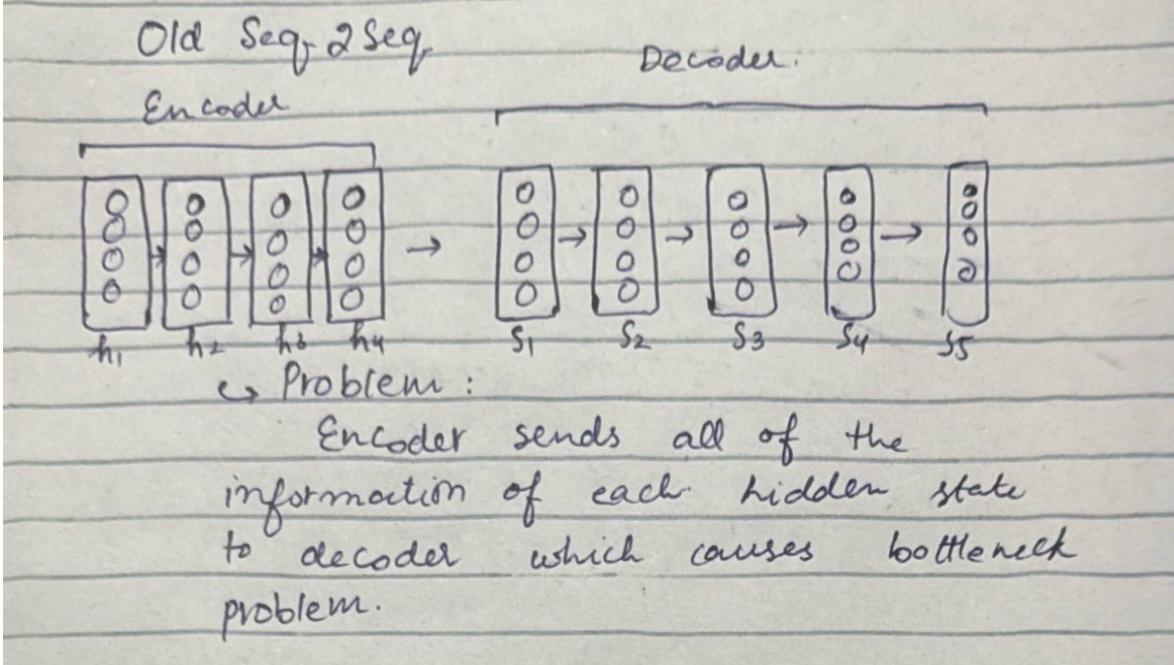
$$P_t(i) = \frac{e^{i/T}}{\sum e^{i/T}}$$

$$P_t(\text{cat}) = \frac{e^{2/0.5}}{e^{2/0.5} + e^{1/0.5} + e^{0.5/0.5}} = 0.84$$

$$P_t(\text{dog}) = \frac{e^{1/0.5}}{e^{2/0.5} + e^{1/0.5} + e^{0.5/0.5}} = 0.11$$

$$P_t(\text{bird}) = \frac{e^{0.5/0.5}}{e^{2/0.5} + e^{1/0.5} + e^{0.5/0.5}} = 0.04$$

Attention in Seq-to-Seg.



Solution? Attention
(focuses on particular parts of encoder)

Seq-to-Seg with Attention

① Take dot product of all hidden layers with the first decoder state to get attention score.

$$e_t = [s_t h_1, \dots]$$

② Take softmax of attention scores to get attention distribution.

$$\alpha_t = \text{softmax}(e_t)$$

③ Take the weighted ~~average~~^{sum} of this attention distribution to get the attention output.

$$a_t = \sum a_t \cdot h_t$$

④ Concatenate the decoder state being used with attention output to get the new encoder state.

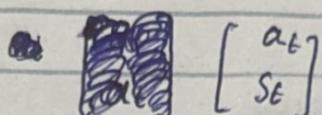


Diagram illustrating the concatenation of a hidden state h_t and an attention score a_t to form a new encoder state SE . The hidden state h_t is represented by a vertical vector of three boxes, and the attention score a_t is represented by a vertical vector of two boxes. The resulting new encoder state SE is shown as a bracketed vertical vector combining the two.

⑤ Do it for all the decoder hidden layers.

Ways of finding attention score.

- ↳ Basic Dot Product
- ↳ Multiplicative
- ↳ Additive.

Benefits

- ↳ improves NMT performance
- ↳ solves bottleneck problem
- ↳ helps with vanishing gradient
- ↳ provides interpretability
- ↳ soft alignment.

* query attends to values.

A way to obtain \rightarrow a fixed size representation of arbitrary set of representations dependent ∇ on some other representations

Example Question:

$$f_1 = [0.1, 0.2, 0.3, 0.4]$$

$$p_2 = [0.2, 0.3, 0.4, 0.5].$$

$$h_3 = [0.3, 0.4, 0.5, 0.6]$$

$$S_t = [0.4, 0.5, 0.6, 0.7]$$

① Attention Scores

$$e_0^t = \text{St.}ht_i$$

$$e_{h_1} = [0.4, 0.5, 0.6, 0.7] [0.1, 0.2, 0.3, 0.4]$$

$$= 0.04 + 0.1 + 0.18 + 0.28$$

0.6

$$e_{H_2} = [0.4, 0.5, 0.6, 0.7] [0.2, 0.3, 0.4, 0.5]$$

0.82.

$$eh3 = [0.4, 0.5, 0.6, 0.7] [0.3, 0.4, 0.5, 0.6]$$

$$\approx 1.04$$

$$e_i^t = [0.6, 0.82, 1.04]$$

② Attention Distribution

$$\alpha^t = \text{softmax}(e^t)$$

$$\alpha_{h1} = \frac{e^{0.6}}{e^{0.6} + e^{0.82} + e^{1.04}} = 0.26$$

$$\alpha_{h2} = \frac{e^{0.82}}{e^{0.6} + e^{0.82} + e^{1.04}} = 0.32$$

$$\alpha_{h3} = \frac{e^{1.04}}{e^{0.6} + e^{0.82} + e^{1.04}} = 0.41$$

$$\alpha^t = [0.26, 0.32, 0.41]$$

③ Attention Output

$$a^t = \sum \alpha_i^t h_i$$

$$a^t = 0.26 \begin{bmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{bmatrix} + 0.32 \begin{bmatrix} 0.2 \\ 0.3 \\ 0.4 \\ 0.5 \end{bmatrix} + 0.41 \begin{bmatrix} 0.3 \\ 0.4 \\ 0.5 \\ 0.6 \end{bmatrix}$$

$$= \begin{bmatrix} 0.026 \\ 0.052 \\ 0.078 \\ 0.104 \end{bmatrix} + \begin{bmatrix} 0.064 \\ 0.096 \\ 0.128 \\ 0.16 \end{bmatrix} + \begin{bmatrix} 0.123 \\ 0.164 \\ 0.205 \\ 0.246 \end{bmatrix}$$

$$= \begin{bmatrix} 0.213 \\ 0.312 \\ 0.411 \\ 0.51 \end{bmatrix}$$

$$a^t = [0.213, 0.312, 0.411, 0.51]$$

④ New Encoder State

$$[a_t][s_t] = [[0.213, 0.312, 0.411, 0.51], [0.4, 0.5, 0.6, 0.7]]$$

Score
 \Rightarrow Attention Calculation Using Multiplicative Attention

Let $W = \begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \\ 0.2 & 0.3 & 0.4 & 0.5 \\ 0.3 & 0.4 & 0.5 & 0.6 \\ 0.4 & 0.5 & 0.6 & 0.7 \end{bmatrix}$

$$e_i = s^T w h_i$$

$$w h_1 = \begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \\ 0.2 & 0.3 & 0.4 & 0.5 \\ 0.3 & 0.4 & 0.5 & 0.6 \\ 0.4 & 0.5 & 0.6 & 0.7 \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{bmatrix} = \begin{bmatrix} 0.3 \\ 0.4 \\ 0.5 \\ 0.6 \end{bmatrix}$$

$$w h_2 = \begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \\ 0.2 & 0.3 & 0.4 & 0.5 \\ 0.3 & 0.4 & 0.5 & 0.6 \\ 0.4 & 0.5 & 0.6 & 0.7 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0.3 \\ 0.4 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.4 \\ 0.5 \\ 0.6 \\ 0.7 \end{bmatrix}$$

$$w h_3 = \begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \\ 0.2 & 0.3 & 0.4 & 0.5 \\ 0.3 & 0.4 & 0.5 & 0.6 \\ 0.4 & 0.5 & 0.6 & 0.7 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0.3 \\ 0.4 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.6 \\ 0.7 \\ 0.8 \end{bmatrix}$$

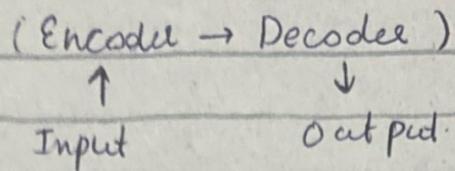
$$e_1 = s_t^T w h_1 \\ = [0.4, 0.5, 0.6, 0.7] \cdot [0.3, 0.4, 0.5, 0.6] \\ = 1.04$$

$$e_2 = 1.26$$

$$e_3 = 1.48$$

$$e^t = [1.04, 1.26, 1.48]$$

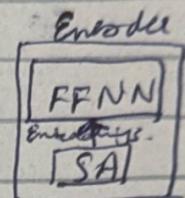
Transformers



↳ Encoder (last hidden state's) 's output is the input of Decoder's (first hidden states)

↳ Sublayers of Encoder

- Feed Forward Neural Network
- Self Attention



↳ Sublayers of Decoder

- Feed Forward NN
- Encoder Decoder Attention
- Self Attention

↳ Path of one encoder layer

Input → SA → Embeddings → FFNN → Output

Question

$$x_1 = [0.1 \ 0.2 \ 0.3 \ 0.4]$$

$$x_2 = [0.5 \ 0.6 \ 0.7 \ 0.8]$$

$$x_3 = [0.9 \ 1.0 \ 1.1 \ 1.2]$$

$$WQ = \begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \\ 0.2 & 0.3 & 0.4 & 0.5 \\ 0.3 & 0.4 & 0.5 & 0.6 \\ 0.4 & 0.5 & 0.6 & 0.7 \end{bmatrix}$$

$$WK = \begin{bmatrix} 0.4 & 0.3 & 0.2 & 0.1 \\ 0.5 & 0.4 & 0.3 & 0.2 \\ 0.6 & 0.5 & 0.4 & 0.3 \\ 0.7 & 0.6 & 0.5 & 0.4 \end{bmatrix}$$

$$W_V = \begin{bmatrix} 0.1 & 0.3 & 0.5 & 0.7 \\ 0.2 & 0.4 & 0.6 & 0.8 \\ 0.3 & 0.5 & 0.7 & 0.9 \\ 0.4 & 0.6 & 0.8 & 1.0 \end{bmatrix}$$

$$\textcircled{1} \quad Q_1 = [0.1 \ 0.2 \ 0.3 \ 0.4] \begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \\ 0.2 & 0.3 & 0.4 & 0.5 \\ 0.3 & 0.4 & 0.5 & 0.6 \\ 0.4 & 0.5 & 0.6 & 0.7 \end{bmatrix} \\ = [0.3 \ 0.4 \ 0.5 \ 0.61]$$

$$K_1 = [0.1 \ 0.2 \ 0.3 \ 0.4] \begin{bmatrix} 0.4 & 0.3 & 0.2 & 0.1 \\ 0.5 & 0.4 & 0.3 & 0.2 \\ 0.6 & 0.5 & 0.4 & 0.3 \\ 0.7 & 0.6 & 0.5 & 0.4 \end{bmatrix}$$

$$= [0.6 \ 0.5 \ 0.4 \ 0.3]$$

$$K_2 V_1 = [0.1 \ 0.2 \ 0.3 \ 0.4] \begin{bmatrix} 0.1 & 0.3 & 0.5 & 0.7 \\ 0.2 & 0.4 & 0.6 & 0.8 \\ 0.3 & 0.5 & 0.7 & 0.9 \\ 0.4 & 0.6 & 0.8 & 1.0 \end{bmatrix} \\ = [0.3 \ 0.5 \ 0.7 \ 0.9]$$

$$K_2 = [0.5 \ 0.6 \ 0.7 \ 0.8] \begin{bmatrix} 0.4 & 0.3 & 0.2 & 0.1 \\ 0.5 & 0.4 & 0.3 & 0.2 \\ 0.6 & 0.5 & 0.4 & 0.3 \\ 0.7 & 0.6 & 0.5 & 0.4 \end{bmatrix} \\ = [1.48 \ 1.22 \ 0.96 \ 0.7]$$

$$K_3 = [0.9 \ 1.0 \ 1.1 \ 1.2] \begin{bmatrix} 0.4 & 0.3 & 0.2 & 0.1 \\ 0.5 & 0.4 & 0.3 & 0.2 \\ 0.6 & 0.5 & 0.4 & 0.3 \\ 0.7 & 0.6 & 0.5 & 0.4 \end{bmatrix} \\ = [2.36 \ 1.94 \ 1.52 \ 1.1]$$

② Scores

$$\begin{aligned}S_1 &= Q_1 \cdot K_1 \\&= [0.3 \ 0.4 \ 0.5 \ 0.6] \cdot [0.6 \ 0.5 \ 0.4 \ 0.3] \\&= 0.18 + 0.2 + 0.2 + 0.18 \\&= 0.76\end{aligned}$$

$$\begin{aligned}S_1 &= Q_1 \cdot K_2 \\&= [0.3 \ 0.4 \ 0.5 \ 0.6] \cdot [1.48 \ 1.22 \ 0.96 \ 0.7] \\&= 1.85\end{aligned}$$

$$\begin{aligned}S_3 &= Q_1 \cdot K_3 \\&= [0.3 \ 0.4 \ 0.5 \ 0.6] \cdot [2.36 \ 1.94 \ 1.52 \ 1.1] \\&= 2.9\end{aligned}$$

③ Div by \sqrt{dk} , Let $\sqrt{dk} = 8$

$$S_{1\text{new}} = 0.76 / 8 = 0.095$$

$$S_{2\text{new}} = 1.85 / 8 = 0.231$$

$$S_{3\text{new}} = 2.9 / 8 = 0.362$$

④ Softmax for $S_{1\text{new}}$

$$\begin{aligned}S_{1\text{new softmax}} &= \frac{e^{0.095}}{e^{0.095} + e^{0.231} + e^{0.362}} \\&= 0.29\end{aligned}$$

⑤ Now calculating z_1

$$z_1 = S_{1, \text{softmax}} \times V_1$$
$$= 0.29 [0.9 \ 0.5 \ 0.7 \ 0.9]$$

$$= [0.087 \ 0.145 \ 0.203 \ 0.261]$$

↳ Multi headed.

Concatenate all z 's and multiply with W_0 to get z .

$$z = z_c \cdot W_0$$

↳ Positional Encoding.

Used for getting diff. probs or results for diff. sentences based on positions.