

SLIDE 1 : RECAP CNNs, RNNs

Slide 2 :

① Diff BETWEEN Seq \rightarrow LABEL
Seq \rightarrow Seq

(Seq \rightarrow LABEL)

Model: $P(y^{(T)} | x^{(1)}, \dots, x^{(T)})$

Seq \rightarrow Seq: $P(y^{(1)}, \dots, y^{(T_y)} | x^{(1)}, \dots, x^{(T_x)})$

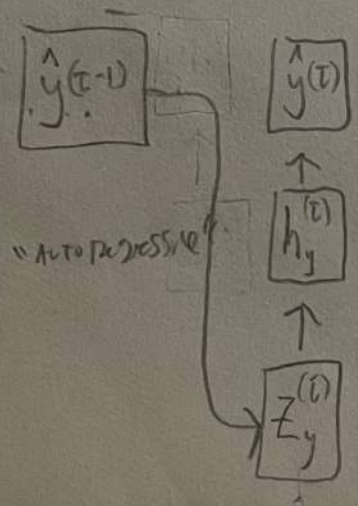
①

EXAMPLES: MACHINE TRANSLATION
 $x^{(1)}, \dots, x^{(T_x)}$ French sentence
 $y^{(1)}, \dots, y^{(T_y)}$ ENGLISH sentence

TEXT GENERATION

$x^{(1)}, \dots, x^{(T)}$ context
 $y^{(1)}, \dots, y^{(T_y)} = x^{(T+1)}, \dots, x^{(T_x)}$

\rightarrow Return to SLIDE



"Auto-regressive"

NOTES: $\hat{y}^{(0)} = \langle \text{START} \rangle$

• DECODER STOPS when $\hat{y} = \langle \text{STOP} \rangle$

THE AUTO REGRESSIVE DECODER

ALGORITHM

① compute $h_x^{(T_x)}(x^{(1)}, \dots, x^{(T_x)})$

② set $T=1$, $\hat{y}^{(0)} = \langle \text{start} \rangle$, $h_y^{(1)} = h_x^{(T_x)}$

③ while true:

$h^{(T)} = f_h(h_y^{(T-1)}, \hat{y}^{(T-1)}) \rightarrow$ compute Hidden State

$y^{(T)} = f_y(h^{(T)}) \rightarrow$ predict T^{th} word

if $y^{(T)} = \langle \text{stop} \rangle$:
Break;

$T = T+1$

NOTE:

SEQ 2 LABEL

$$h^{(T)} = f_h(h^{(T-1)}, x^{(T)})$$

comes
from
the
data

Seq 2 Seq

$$h^{(T)} = f_h(h^{(T-1)}, y^{(T-1)})$$

comes
from
the
model

THE "INFERENCE GAP"

TRAINING TASK: $\hat{\theta} = \underset{\theta}{\text{argmax}} P(y^{(T)} | f_h(h^{(T-1)}, y^{(T-1)}); \theta)$

PREDICTION TASK: $\hat{y}^{(1)}, \dots, \hat{y}^{(T)} = \underset{y^{(1)}, \dots, y^{(T)}}{\text{argmax}} P(y^{(1)}, \dots, y^{(T)} | f_h(h_x^{(T_x)}, y^{(0)}); \theta)$

OUR TRAINING AND PREDICTION TASKS have DIVERGED!

TEACHER FORCING :

$$L_{CE} (y^{(t)}, p(y^{(t)} | y^{(n)}, \dots, y^{(t-1)}))$$

$y^{(t)} = t^{th}$ word in TRAINING EXAMPLE

$\hat{y}^{(t)} = t^{th}$ prediction

\Rightarrow TRAINED ONLY ON GROUND TRUTH CONTEXT WORDS

SCHEDULED Sampling :

$$L_{CE} (y^{(t)}, p(y^{(t)} | \tilde{y}^{(n)}, \dots, \tilde{y}^{(t-1)}))$$

where
$$\tilde{y}^{(t)} = \begin{cases} \hat{y}^{(t)} & \text{with probability } p \\ y^{(t)} & \text{with prob } 1-p \end{cases}$$

\hookrightarrow EVEN WITH SCHEDULED Sampling: THE TRAINING AND INTERLUCE DISTRIBUTIONS ARE DIFFERENT, UNLESS θ IS A PERFECT MODEL OF THE DATA.

CAN WE JUST ESTIMATE $p(y^{(n)}, \dots, y^{(t)})$ USING MLE?

(4)

$$p(y^{(n)}, \dots, y^{(t_y)}) = \frac{1}{T_y} \prod_{t=1}^{T_y} p(y^{(t)} | y^{(n)}, \dots, y^{(t-1)}; \theta)$$

\hookrightarrow No, SCALES EXPONENTIALLY WITH SEQ LEN T_y

\hookrightarrow EX: $N=100,000$ Vocab SIZE $T_y=5 \Rightarrow$ W'D NEED TO EVAL $(10^5)^5 = 10^{25}$ POSSIBILITIES!

INFERENCE GAP Cont...

L9.

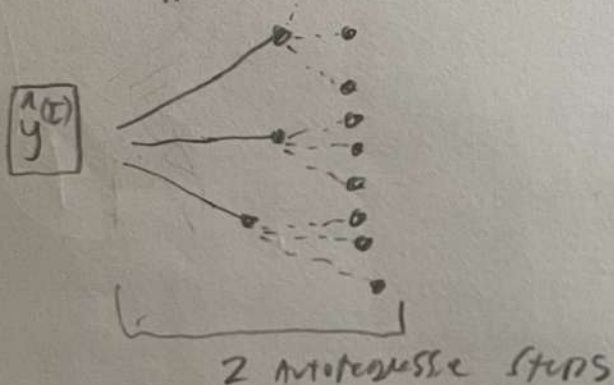
APPROXIMATE SEARCH ALGORITHMS:

5

BEAM SEARCH:

IDEA: AT EVERY AUTO-REGRESSIVE STEP, ONLY
CONSIDER TOP-K MOST PROBABLE WORDS

EX; $K=3$
 $N=10^5$



ORIGINAL PROBLEM:

$$N^2 = 10^{10}$$

$$K^2 = 9$$

IN PRACTICE THIS WORKS FAIRLY WELL.

2 PRIMARY CHALLENGES WITH SEQ2SEQ MODELS:

- ① THE INFORMATION "BOTTLENECK"
- ② THE INFERENCE GAP

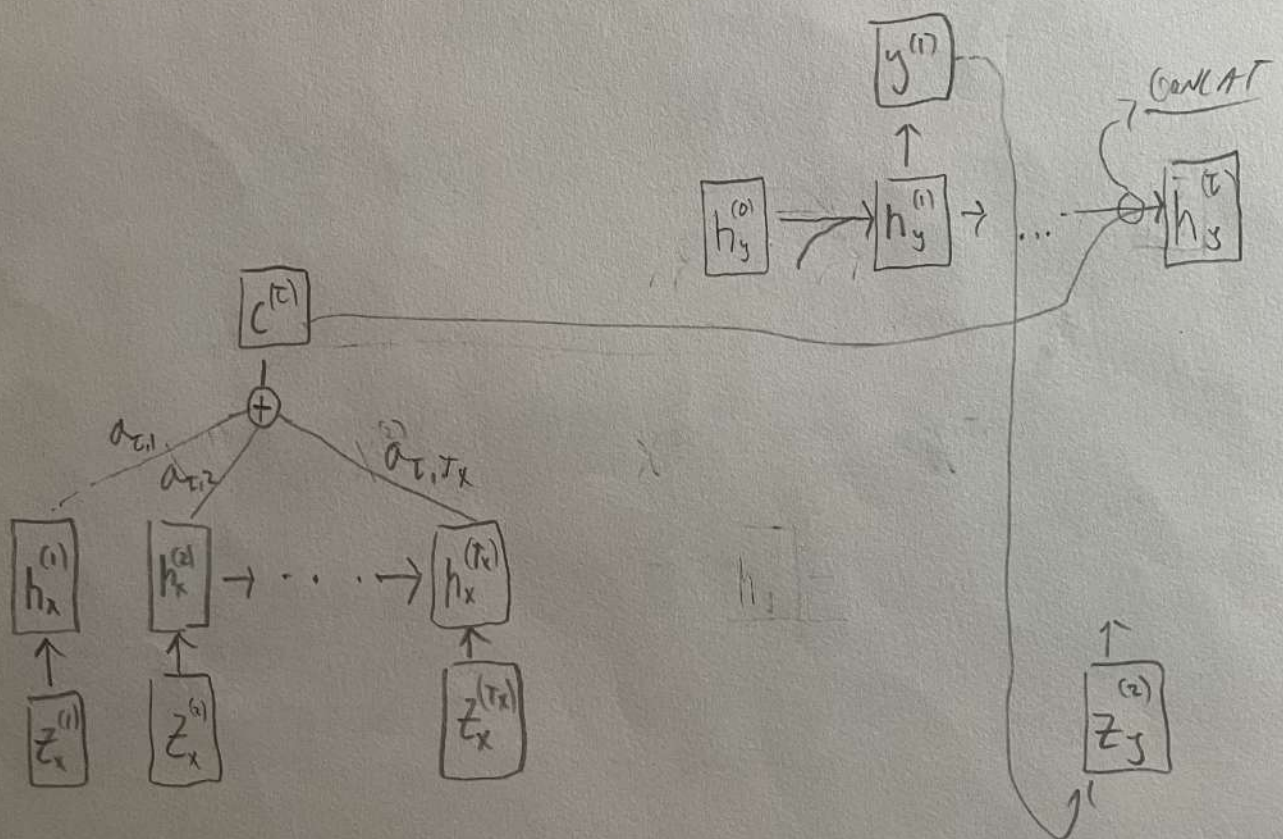
BOTH OF THESE CHALLENGES HAVE BEEN ADDRESSED
WITH THE ADVENT OF "ATTENTION"!

FIRST: THE INFORMATION BOTTLENECK

ATTENTION BASED SEQ2SEQ MODELING (2014) | L9-5

→ STANDARD BIDIRECTIONAL LSTM BASED E-D ARCH WITH ONE KEY ADDITION: AN ATTENTION MECHANISM BETWEEN THE DECODER AND EACH HIDDEN STATE IN THE ORIGINAL SEQ X.

$$h^{(t)} = \{ \vec{h}^{(t)}, \overleftarrow{h}^{(t)} \}$$



Attention BASED Seq2Seq Models Cont.

So total S $a_{t,t'}$?

$$A \in [0,1]^{T_y \times T_x}$$

(2)

$$\sum_{t'=1}^{T_x} a_{t,t'} = 1$$

$$(1) c^{(t)} = \sum_{t'=1}^{T_x} a_{t,t'} h_x^{(t')}$$

$\rightarrow a_{t,t'}^{(t,t')} = \text{DEGREE of ATTENTION } y^{(t)} \text{ has on } x^{(t')}$

$$a_{t,t'} = \frac{e^{\alpha_{t,t'}}}{\sum_{t''=1}^{T_x} e^{\alpha_{t,t''}}}$$

$\alpha =$ (1) DOT-PRODUCT $\alpha_{t,t'} = h_y^{(t)} \cdot h_x^{(t')}$

(2) NEURAL NET $\alpha_{t,t'} = f_{h_y}(h_y^{(t)}) \cdot f_{h_x}(h_x^{(t')})$

\hookrightarrow Where $f_{\theta}(\cdot; \theta) = \text{NN}$

Attention is all you need cont. -

L9-9

→ Notice that each input $h^{(i)}$ is used in 3 different ways!

② → Attention as a Fuzzy lookup table
for every $c^{(i)}$ there is

① 1 Query: $h^{(i)}$

② T separate keys: $h^{(1)}, \dots, h^{(T)}$

③ T separate values: $h^{(1)}, \dots, h^{(T)}$

$$\begin{aligned} \rightarrow c^{(i)} &= \text{Dict}[\text{Query}^{(i)}] \text{ \& fuzzy lookup} \\ &= \sum_{i'=1}^T \text{Value}^{(i')} \cdot \sigma_{\text{softmax}}(\text{Query}^{(i)}, \text{Key}^{(i')}) \\ &= \sum_{i'=1}^T h^{(i)} \cdot \sigma_{\text{softmax}}(h^{(i)}, h^{(i')}) \end{aligned}$$

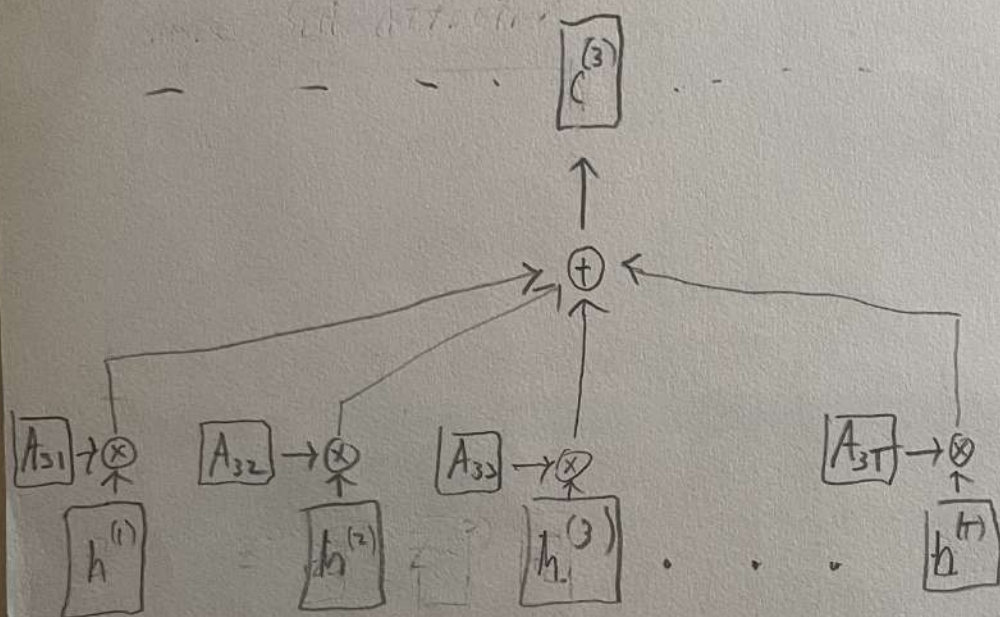
↳ The context vector at t^{th} step
is the sum of all inputs $h^{(1)} \dots h^{(T)}$
weighted by the similarity factor $\sigma_{\text{softmax}}(h^{(i)}, h^{(i')})$

Simple Self Attention

L9.7

Slide 7

IDEA: Could we just Rely ON ATTENTION
ALONE AND DO AWAY WITH RNNs?



$$c^{(T)} = \sum_{i'} A_{T,i'} h^{(i')} \quad \text{where} \quad A_{T,i'} = \frac{e^{h^{(T)} \cdot h^{(i')}}}{\sum_{i''} e^{h^{(T)} \cdot h^{(i'')}}}$$

Where Sum over any
column of $A = 1$

Simple Self Attention Cont...

NOTE: Simple Self Attention is NOT a Sequence Model (NECESSARILY), depends ON ~~whether~~ how it is computed, the Simple Self Attention Model is a Set Model.

ATTENTION IS ALL YOU NEED

~~THE TRANSFORMER MODEL~~

INTRODUCED THE TRANSFORMER MODEL (Slide 8)

How'd They get it to work?

① Scaled Dot-Product Self Attention

↳ Problem: We're computing a Soft MAX over DOT PRODUCTS: the magnitude of the $\langle \cdot \rangle$ scales with \sqrt{D}

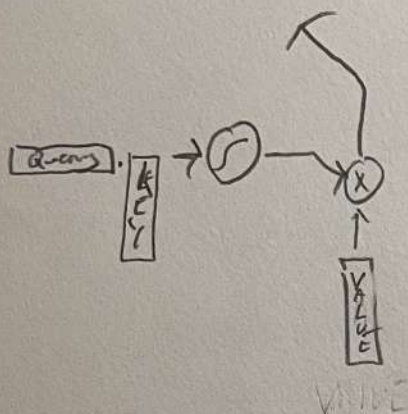
↳ ex:

$$\begin{aligned}\sum [1, 1]^2 &= \sqrt{2} \\ \sum [1, 1, 1]^2 &= \sqrt{3} \\ \sum [1, 1, 1, 1] &= \sqrt{4}\end{aligned}$$

So: $A \rightarrow \frac{A}{\sqrt{D}}$

The K, Q, V I'd like more generality:

(C9)



Each input h is serves
as the key, query, and
value

③ ~~the previous set of parameters~~ A Learnable Set of
keys, queries, values

$$K_T = K h^{(i)} + b_K$$

$$Q_T = Q h^{(i)} + b_Q$$

$$V_T = V h^{(i)} + b_V$$

$$\Theta = \{K, Q, V, b_K, b_Q, b_V\}$$

④ Multi-head SA

$$K = \{K^{(1)}, \dots, K^{(n)}\}$$

$$Q = \{Q^{(1)}, \dots, Q^{(n)}\}$$

$$V = \{V^{(1)}, \dots, V^{(n)}\}$$

↳ this gives us n context vectors at each
sequence position i : $c^{(i)} = [c_1^{(i)}, \dots, c_n^{(i)}]$

↳ in practice we can make each
 $K, Q, V \in \mathbb{R}^{d/n}$ such that $c^{(i)} \in \mathbb{R}^d$