

DS-GA 1011 NLP

Fall 2023

Recitation 4

Encoder-Decoder Model for Machine Translation

Lavender Jiang

Sep 29, 2023

## German

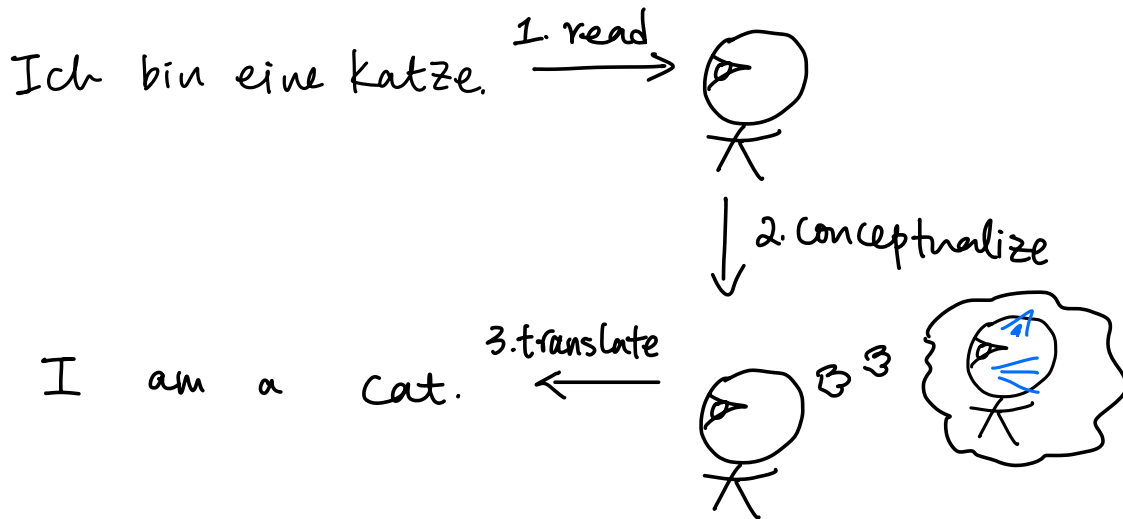
Ich bin eine Katze.  
Ich bin eine Frau.  
Ich bin Bürger.

## English

I am a cat.  
I am a woman.  
I am a citizen.

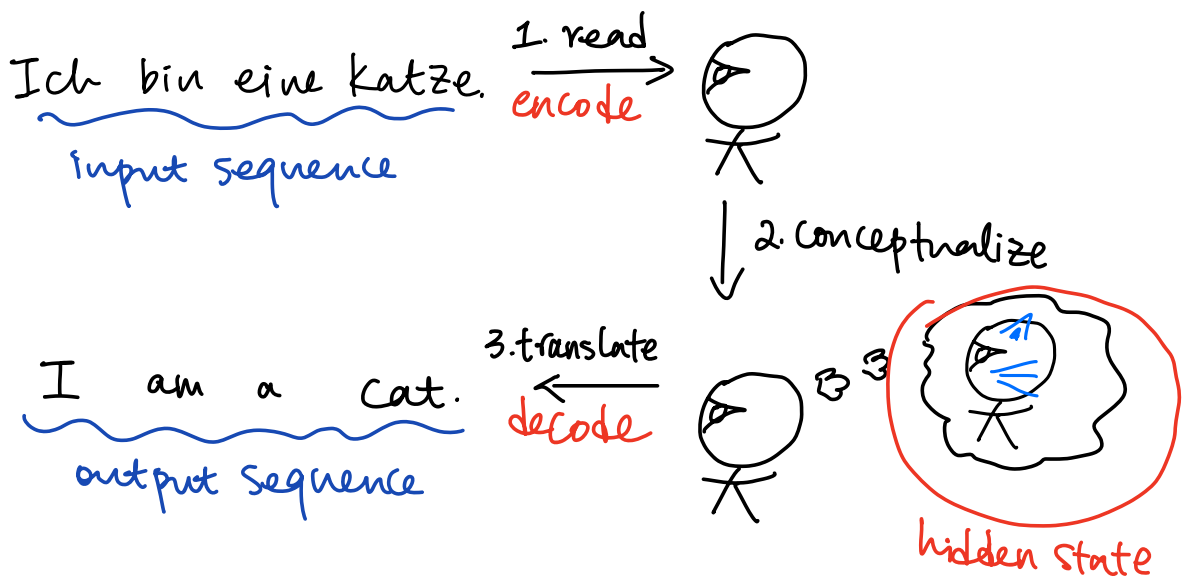
1

---



2

---



3

Attention

Ich bin eine Katze.  
input sequence

1. read  
encode



2. conceptualize



hidden state

I am a Cat.  
output sequence

3. translate  
decode

4

Attention

Ich bin eine Katze.  
input sequence

1. read  
encode



2. conceptualize




hidden state


I am a Cat.  
output sequence

3. translate  
decode

5

Ich bin eine Katze.  $\xrightarrow{1. \text{ read}}$  

6

Ich bin eine Katze.  $\xrightarrow[1. \text{ read}]{\text{tokenize + parameterized encoder}}$  

7

Ich bin eine Katze  
 $\downarrow$  tokenize  
 [1, 3, 4, 5]

vocab	id
Ich	1
du	2
bin	3
eine	4
Katze	5
frau	6

8

<bos> Ich bin eine Katze <eos>  
 $\downarrow$  tokenize  
 [7, 1, 3, 4, 5, 8]

vocab	id
Ich	1
du	2
bin	3
eine	4
Katze	5
frau	6
<bos>	7
<eos>	8

special tokens {

9

<bos> Ich bin eine Katze <eos>

↓ tokenize

[7, 1, 3, 4, 5, 8]

<bos> Ich bin eine Studentin <eos>

↓ tokenize

[7, 1, 3, 4, 9, 10, 8]

Vocab	id
Ich	1
du	2
bin	3
eine	4
Katze	5
frau	6
<bos>	7
<eos>	8
student	9
##in	10

10

Suppose encoder is linear. i.e.  $enc := A$ .

vocab size = 10, hidden size = 20.  $A \in \mathbb{R}^{10 \times 20}$

Then hidden for input  $s$  is  $h = enc(\text{tokenize}(s))$   
 $= enc(x) = Ax$

<bos> Ich bin eine Katze <eos>  $s$

↓ tokenize

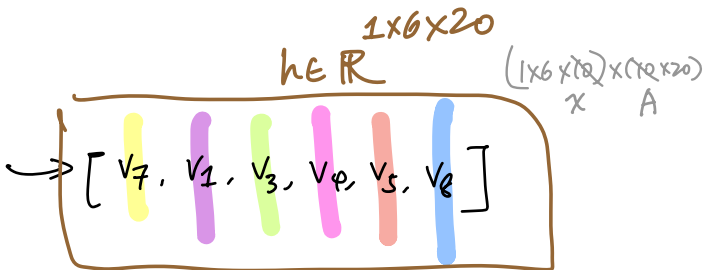
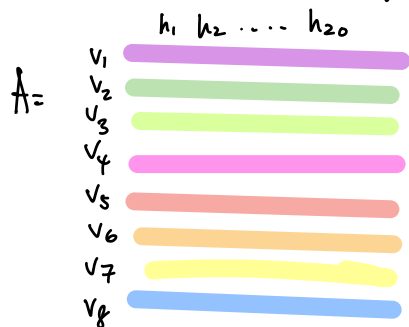
[7, 1, 3, 4, 5, 8]  $x \in \mathbb{R}^{1 \times 6}$

each token has a learned embedding

index for embedding

$x' \in \mathbb{R}^{1 \times 6 \times 10}$

$x' = [e_7, e_1, e_3, e_4, e_5, e_8]$



11

$\langle \text{bos} \rangle$  Ich bin eine katze  $\langle \text{eos} \rangle$   $S_1$

↓ tokenize

[7, 1, 3, 4, 5, 8]  $X \in \mathbb{R}^{1 \times 6}$

$\langle \text{bos} \rangle$  Ich bin eine studentin  $\langle \text{eos} \rangle$   $S_2$

↓ tokenize

[7, 1, 3, 4, 9, 10, 8]  $X \in \mathbb{R}^{1 \times 7}$

$X = ?$

$X \in \mathbb{R}^{B \times M}$

B: batch size  
M: max length

12

$M=8, B=2$

[ $\langle \text{bos} \rangle$  Ich bin eine katze  $\langle \text{eos} \rangle$   $\langle \text{pad} \rangle$ ]

[ $\langle \text{bos} \rangle$  Ich bin eine studentin  $\langle \text{eos} \rangle$ ]

$M=9, B=2$

[ $\langle \text{bos} \rangle$  Ich bin eine katze  $\langle \text{eos} \rangle$   $\langle \text{pad} \rangle$   $\langle \text{pad} \rangle$ ]

[ $\langle \text{bos} \rangle$  Ich bin eine studentin  $\langle \text{eos} \rangle$   $\langle \text{pad} \rangle$ ]

13

batch of size 2 { [ $\langle \text{bos} \rangle$  Ich bin eine katze  $\langle \text{eos} \rangle$   $\langle \text{pad} \rangle$   $\langle \text{pad} \rangle$ ]  
[ $\langle \text{bos} \rangle$  Ich bin eine studentin  $\langle \text{eos} \rangle$   $\langle \text{pad} \rangle$ ]

↓  
tokenizer

↓  
encoder

↓  
hidden

↓  
decoder

[ $\langle \text{bos} \rangle$  I am]  
[ $\langle \text{bos} \rangle$  I am]

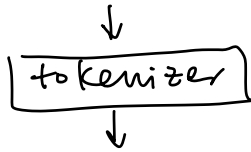
[ $\langle \text{bos} \rangle$ ]  
[ $\langle \text{bos} \rangle$ ]

[ $\langle \text{bos} \rangle$  I]  
[ $\langle \text{bos} \rangle$  I]

How can you increase batch size under the same compute?

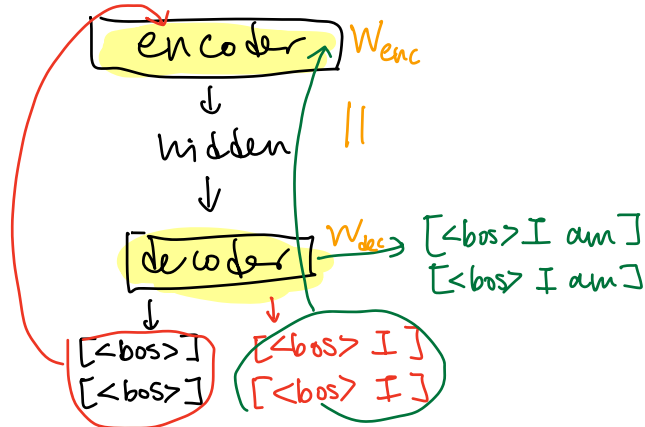
14

batch of size 2 { [ $\langle \text{bos} \rangle$  Ich bin eine katze  $\langle \text{eos} \rangle$   $\langle \text{pad} \rangle$   $\langle \text{pad} \rangle$  ]  
 [ $\langle \text{bos} \rangle$  Ich bin eine student ##in  $\langle \text{eos} \rangle$   $\langle \text{pad} \rangle$  ]



$$W^{t+1} = \eta (\nabla W_{\text{enc}} + \nabla W_{\text{dec}})$$

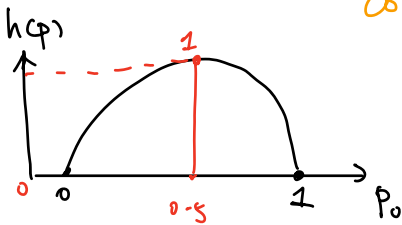
weight sharing!  
 also helps w  
 • overfitting  
 • robustness  
 • semantic matching



14

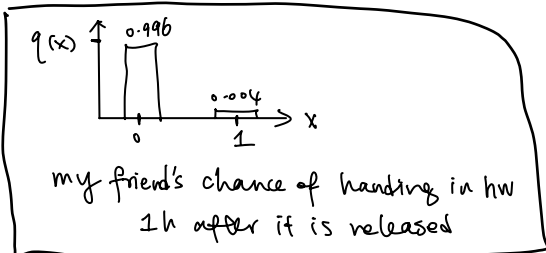
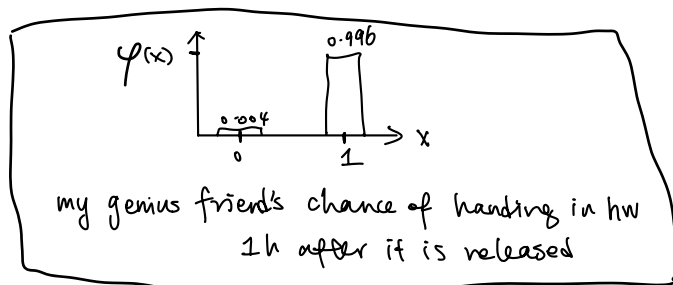
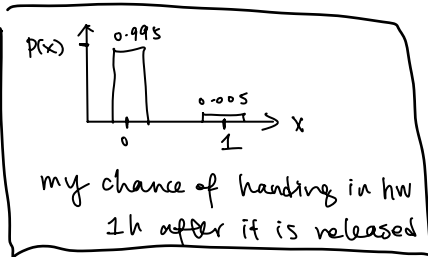
entropy measures the "uncertainty" / "chaoticness" of a dist.

coin flip  $\$1$



$$h(p) = - \sum_i p_i \log p_i = \mathbb{E}_p [-\log p] \in [0, 1]$$

Cross entropy measures how "close" two distributions are



$$H(p, q) = - \sum_i p_i \log q_i = \mathbb{E}_p [-\log q_i]$$

$$KL(p, q) = - \sum_i p_i \log \left( \frac{q_i}{p_i} \right) = \mathbb{E}_p \left[ -\log \frac{q_i}{p_i} \right]$$

$$H(p, q) = KL(p, q) - H(p)$$

$$\text{minimize } H(p, q) \Rightarrow \text{minimize } KL(p, q)$$

15

In our case,  $p = P(y_{t+1} | \underbrace{\text{context}}_{y_{\leq t}, \vec{x}})$ ,

$$q = P(\hat{y}_{t+1} | \text{context}_t)$$

Concrete example:  $p = P(y_{t+1} | [I], \text{Lich, bin, eine, Katz})$



$q = P(\hat{y}_{t+1} | [I], \text{Lich, bin, eine, Katz})$

