

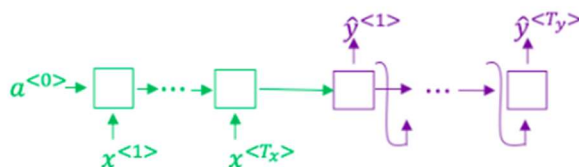
## Your grade: 90%

Your latest: 90% • Your highest: 90% • To pass you need at least 80%. We keep your highest score.

Next item →

1. Consider using this encoder-decoder model for machine translation.

1 / 1 point



True/False: This model is a “conditional language model” in the sense that the encoder portion (shown in green) is modeling the probability of the input sentence  $x$ .

- ☐ True  
☒ False

✓ Correct

The encoder-decoder model for machine translation models the probability of the output sentence  $y$  conditioned on the input sentence  $x$ . The encoder portion is shown in green, while the decoder portion is shown in purple.

2. In beam search, if you decrease the beam width  $B$ , which of the following would you expect to be true? Select all that apply.

☒ Beam search will generally find better solutions (i.e. do a better job maximizing  $P(y|x)$ ).

✗ This should not be selected

As the beam width decreases, beam search runs more quickly, uses up less memory, and converges after fewer steps, but will generally not find the maximum  $P(y|x)$ .

☒ Beam search will use up more memory.

✗ This should not be selected

As the beam width decreases, beam search runs more quickly, uses up less memory, and converges after fewer steps, but will generally not find the maximum  $P(y|x)$ .

☐ Beam search will run more quickly.

☐ Beam search will converge after fewer steps.

(I MISREAD “increase the beam width”)

3. In machine translation, if we carry out beam search without using sentence normalization, the algorithm will tend to output overly short translations.

- ☐ False  
☒ True

✓ Correct

4. Suppose you are building a speech recognition system, which uses an RNN model to map from audio clip  $x$  to a text transcript  $y$ . Your algorithm uses beam search to try to find the value of  $y$  that maximizes  $P(y | x)$ .

On a dev set example, given an input audio clip, your algorithm outputs the transcript  $\hat{y}$  = "I'm building an A Eye system in Silly con Valley.", whereas a human gives a much superior transcript  $y^*$  = "I'm building an AI system in Silicon Valley."

According to your model,

$$P(\hat{y} | x) = 7.21 \times 10^{-8}$$

$$P(y^* | x) = 1.09 \times 10^{-7}$$

Would you expect increasing the beam width  $B$  to help correct this example?

- ☒ Yes, because  $P(y^* | x) > P(\hat{y} | x)$  indicates the error should be attributed to the search algorithm rather than to the RNN.
- ☐ Yes, because  $P(y^* | x) > P(\hat{y} | x)$  indicates the error should be attributed to the RNN rather than to the search algorithm.
- ☐ No, because  $P(y^* | x) > P(\hat{y} | x)$  indicates the error should be attributed to the RNN rather than to the search algorithm.
- ☐ No, because  $P(y^* | x) > P(\hat{y} | x)$  indicates the error should be attributed to the search algorithm rather than the RNN.

✓ Correct

$P(y^* | x) > P(\hat{y} | x)$  indicates the error should be attributed to the search algorithm rather than to the RNN. Increasing the beam width will generally allow beam search to find better solutions.

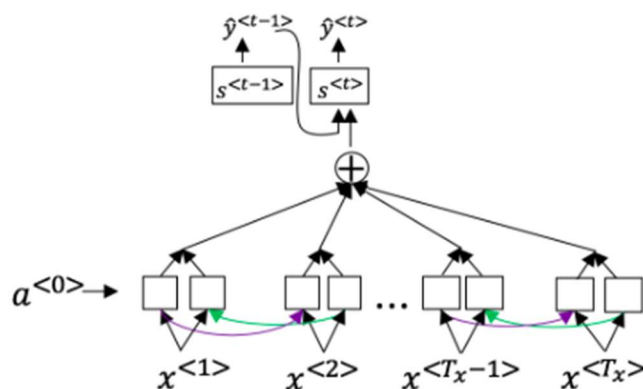
5. Continuing the example from Q4, suppose you work on your algorithm for a few more weeks, and now find that for the vast majority of examples on which your algorithm makes a mistake,  $P(y^* | x) > P(\hat{y} | x)$ . This suggests you should focus your attention on improving the RNN.

- ☐ True  
☒ False

✓ Correct

$P(y^* | x) > P(\hat{y} | x)$  indicates the error should be attributed to the search algorithm rather than to the RNN.

6. Consider the attention model for machine translation.



Further, here is the formula for  $\alpha^{<t,t'>}$ .

$$\alpha^{<t,t'>} = \frac{\exp(e^{<t,t'>})}{\sum_{t'=1}^{T_x} \exp(e^{<t,t'>})}$$

Which of the following statements about  $\alpha^{<t,t'>}$  are true? Check all that apply.

☒  $\sum_{t'} \alpha^{<t,t'>} = 1$  (Note the summation is over  $t'$ .)

☒ Correct

Correct! If we sum over  $\alpha^{<t,t'>}$  for all  $t'$  (the formulation can be seen in the image), the numerator will be equal to the denominator, therefore,  $\sum_{t'} \alpha^{<t,t'>} = 1$ .

☐  $\sum_t \alpha^{<t,t'>} = 1$  (Note the summation is over  $t$ .)

☒ We expect  $\alpha^{<t,t'>}$  to be generally larger for values of  $a^{<t'>}$  that are highly relevant to the value the network should output for  $y^{<t>}$ . (Note the indices in the superscripts.)

☒ Correct

Correct!  $\alpha^{<t,t'>}$  is equal to the amount of attention  $y^{<t>}$  should pay to  $a^{<t'>}$ . So, if a value of  $a^{<t'>}$  is highly relevant to  $y^{<t>}$ , then the attention coefficient  $\alpha^{<t,t'>}$  should be larger. Note the difference between  $a$  (activation) and  $\alpha$  (attention coefficient).

☐ We expect  $\alpha^{<t,t'>}$  to be generally larger for values of  $a^{<t>}$  that are highly relevant to the value the network should output for  $y^{<t'>}$ . (Note the indices in the superscripts.)

7. The network learns where to “pay attention” by learning the values  $e^{<t,t'>}$ , which are computed using a small neural network:

We can't replace  $s^{<t-1>}$  with  $s^{<t>}$  as an input to this neural network. This is because  $s^{<t>}$  depends on  $\alpha^{<t,t'>}$  which in turn depends on  $e^{<t,t'>}$ ; so at the time we need to evaluate this network, we haven't computed  $s^{<t>}$  yet.

☒ True

☐ False

☒ Correct

8. The attention model performs the same as the encoder-decoder model, no matter the sentence length.

☒ False

☐ True

✓ **Correct**

The performance of the encoder-decoder model declines as the amount of words increases. The attention model has the greatest advantage when the input sequence length  $T_x$  is large.

9. Under the CTC model, identical repeated characters not separated by the "blank" character (`_`) are collapsed. Under the CTC model, what does the following string collapse to?

kk\_eee\_\_\_\_ee\_p\_\_eeeeeeee\_\_\_\_rrrrr

☐ kkeeeeepeeeeeeeerrrrr

☐ keper

☒ keeper

☐ ke epe r

✓ **Correct**

The basic rule for the CTC cost function is to collapse repeated characters not separated by "blank". If a character is repeated, but separated by a "blank", it is included in the string.

10. In trigger word detection,  $x^{<t>}$  is:

☐ Whether the trigger word is being said at time  $t$ .

☒ Features of the audio (such as spectrogram features) at time  $t$ .

☐ Whether someone has just finished saying the trigger word at time  $t$ .

☐ The  $t$ -th input word, represented as either a one-hot vector or a word embedding.

✓ **Correct**