

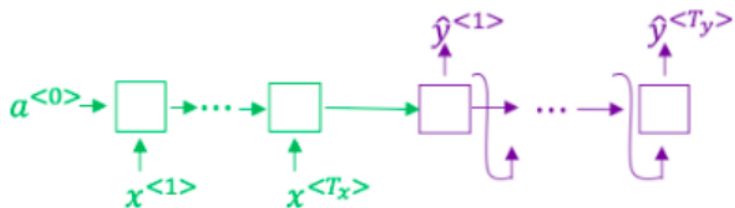
✔ Congratulations! You passed!

[Go to next item](#)

Grade received 90% Latest Submission Grade 90% To pass 80% or higher

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 decoder portion (shown in green) is modeling the probability of the input sentence x .

☐ True

☒ False

 **Expand**

 **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 run more quickly.

✓ **Correct**

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 converge after fewer steps.

☒ 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 $P(y|x)$.

 **Expand**

 **Incorrect**

You didn't select all the correct answers

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

1 / 1 point

☐ False

☒ True

 **Expand**

 **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 $P(y \mid x)$.

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,

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

$$P(y^* \mid x) = 7.21 * 10^{-8}$$

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

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



Expand



Correct

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, $P(y^* \mid x) > P(\hat{y} \mid x)$. This suggests you should focus your attention on improving the search algorithm.

☐ False.

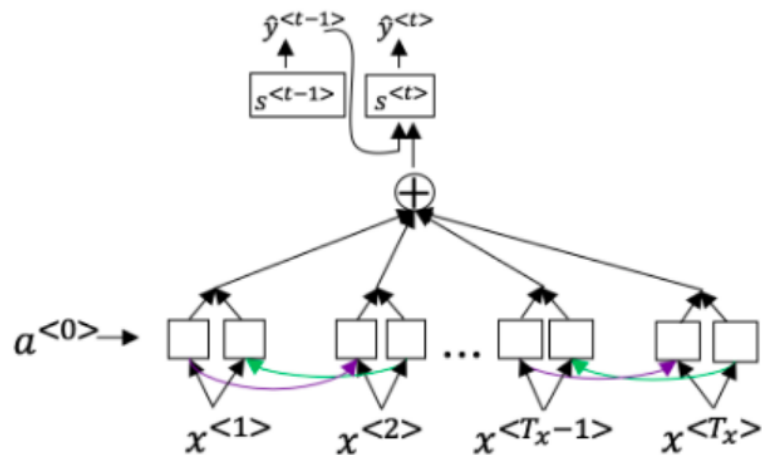
☒ True.

 Expand

 Correct

6. Consider the attention model for machine translation.

1 / 1 point



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.

☐ $\sum_{t'} \alpha^{<t,t'>} = -1$

☐ $\sum_{t'} \alpha^{<t,t'>} = 0$

☒ $\alpha^{<t,t'>}$ is equal to the amount of attention $y^{<t>}$ should pay to $a^{<t'>}$

✓ Correct

Correct! $\alpha^{<t,t'>} =$ amount of attention $y^{<t>}$ should pay to $a^{<t'>}$

☐ 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.)

 **Expand**



Correct

Great, you got all the right answers.

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 replace $s^{<t-1>}$ with $s^{<t>}$ as an input to this neural network because $s^{<t>}$ is independent of $\alpha^{<t,t'>}$ and $e^{<t,t'>}$.

☐ True

☒ False

 Expand

 **Correct**

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 and $e^{<t,t'>}$; so at the time we need to evaluate this network, we haven't computed $s^{<t>}$.

8. Compared to the encoder-decoder model shown in Question 1 of this quiz (which does not use an attention mechanism), we expect the attention model to have the least advantage when:

- ☐ The input sequence length T_x is large.
- ☒ The input sequence length T_x is small.

 Expand

 Correct

The encoder-decoder model works quite well with short sentences. The true advantage for the attention model occurs when the input sentence 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?

__c_oo_o_kk__b_ooooo__oo__kkk

- ☐ cokbok
- ☐ cook book
- ☒ cookbook
- ☐ coookkboooooookkk

 **Expand**

 **Correct**

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

- ☒ Features of the audio (such as spectrogram features) at time t .
- ☐ Whether the trigger word is being said 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.

 Expand

 Correct