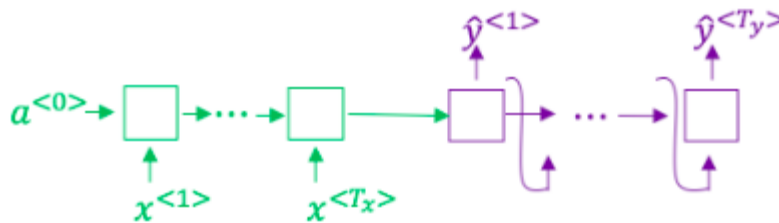


Sequence Models & Attention Mechanism

1.

Question 1

Consider using this encoder-decoder model for machine translation.



True/False: This model is a “conditional language model” in the sense that the decoder portion (shown in purple) is modeling the probability of the output sentence y given the input sentence x .

☐ False

☒ True

[Expand](#)

✓ **Correct**

The encoder-decoder model for machine translation models the probability of the output sentence y conditioned on the input sentence x .

2.

Question 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 converge after fewer steps.

✓ **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 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.

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

[↗ Expand](#)

✓ **Correct**

Great, you got all the right answers.

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

1 / 1 point

☒ False

☐ True

[↗ Expand](#)

✓ **Correct**

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

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 $P(y | x)$.

1 / 1 point

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,

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

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

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

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

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

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

1 / 1 point

- ☒ True.
- ☐ False.

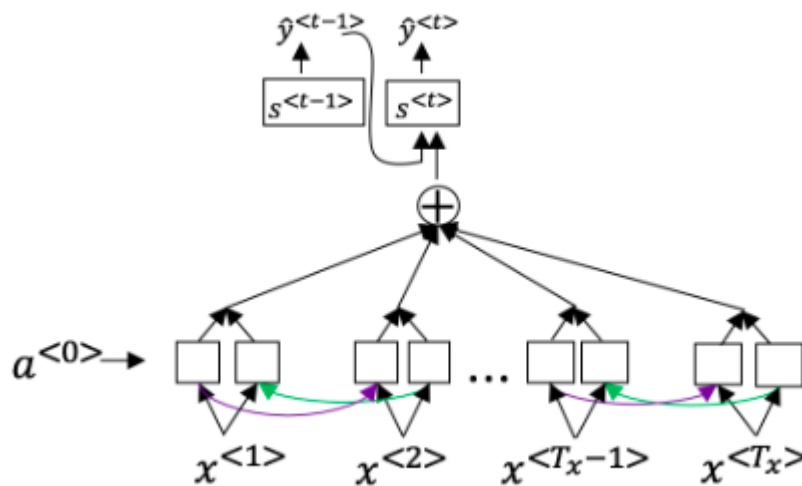
Expand

 **Correct**

6.

Question 6

Consider the attention model for machine translation.



Further, here is the formula for $\langle \cdot, \cdot \rangle_{\alpha_{\langle t, t' \rangle}}$.

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

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

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

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

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

! This should not be selected

We expect $\alpha^{<t,t'>}$ to be generally larger for values of $\alpha^{<t'>}$ that are highly relevant to the value the network should output for $y^{<t>}$, not for $y^{<t'>}$

- ☒ $\sum_{t'} \alpha^{<t,t'>} = 0$. (Note the summation is over t' .)

! This should not be selected

We expect $\alpha^{<t,t'>}$ to be larger for activation values that are highly relevant to the value the network should output for $y^{<t>}$.

the network should output for $y^{<t>}$.

- ☐ $\sum_{t'} \alpha^{<t,t'>} = 1$. (Note the summation is over t' .)
- ☐ $\alpha^{<t,t'>}$ is equal to the amount of attention $y^{<t>}$ should pay to $\alpha^{<t'>}$

 Expand

⊗ Incorrect

You didn't select all the correct answers

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

0 / 1 point

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](#)

✗ Incorrect

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 greatest advantage when:

1 / 1 point

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

[Expand](#)

✓ Correct

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

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

[Expand](#)

✓ Correct

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

1 / 1 point

- ☐ 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.

[Expand](#)

✓ Correct