

Quiz for RNN & Sequence Models

Thursday, May 9, 2019

12:39 AM

1. Question 1

Suppose your training examples are sentences (sequences of words). Which of the following refers to the j^{th} word in the i^{th} training example?



$x^{(i)<j>}x^{(i)<j>}$



$x^{<i>(j)}x^{<i>(j)}$



$x^{(j)<i>}x^{(j)<i>}$



$x^{<j>(i)}x^{<j>(i)}$

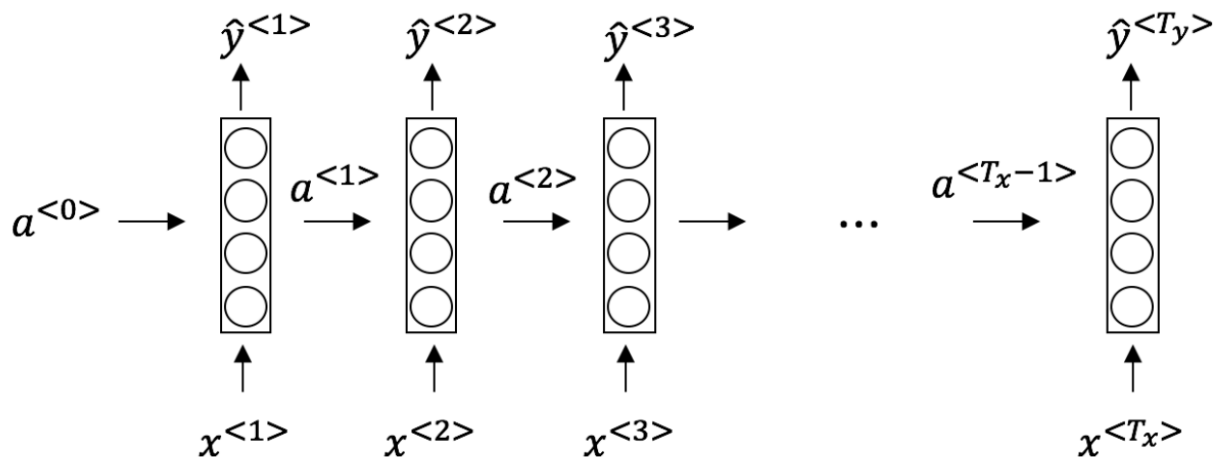
Question 2

1

point

2. Question 2

Consider this RNN:



This specific type of architecture is appropriate when:



$T_x = T_y$



$T_x < T_y$



$T_x > T_y$



$T_x = 1$

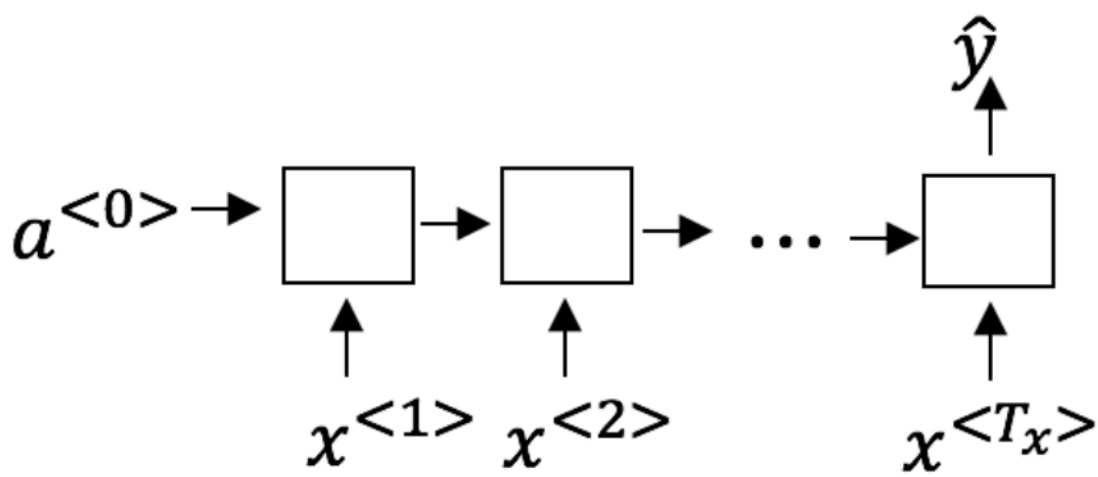
Question 3

1

point

3. Question 3

To which of these tasks would you apply a many-to-one RNN architecture? (Check all that apply).



- ☐ Speech recognition (input an audio clip and output a transcript)
- ☒ Sentiment classification (input a piece of text and output a 0/1 to denote positive or negative sentiment)
- ☐ Image classification (input an image and output a label)
- ☒ Gender recognition from speech (input an audio clip and output a label indicating the speaker’s gender)

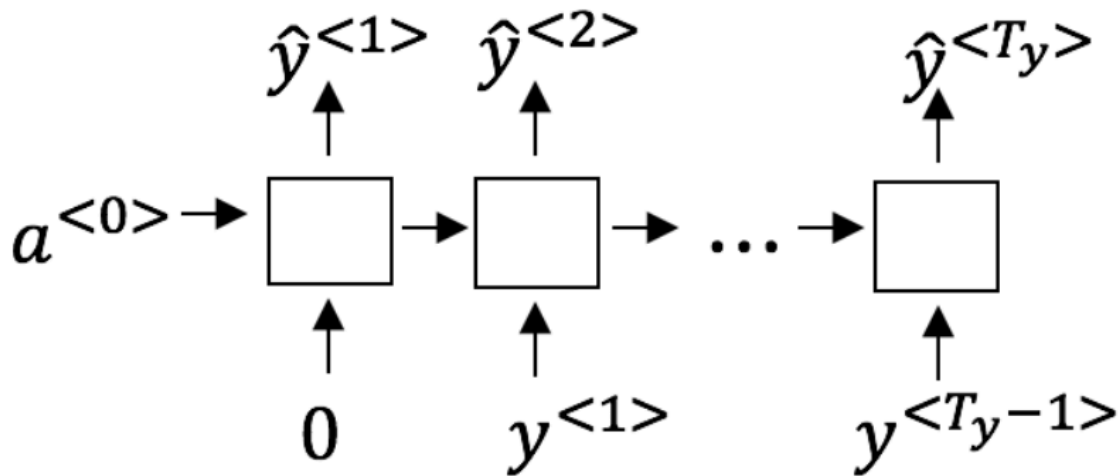
Question 4

1

point

4. Question 4

You are training this RNN language model.



At the t^{th} time step, what is the RNN doing? Choose the best answer.

- ☐ Estimating $P(y_{<1>}, y_{<2>}, \dots, y_{<t-1>})$
- ☐ Estimating $P(y^{\{<t>\}})P(y_{<t>})$
- ☒ Estimating $P(y_{<t>} | y_{<1>}, y_{<2>}, \dots, y_{<t-1>})$
- ☐ Estimating $P(y_{<t>} | y_{<1>}, y_{<2>}, \dots, y_{<t>})$

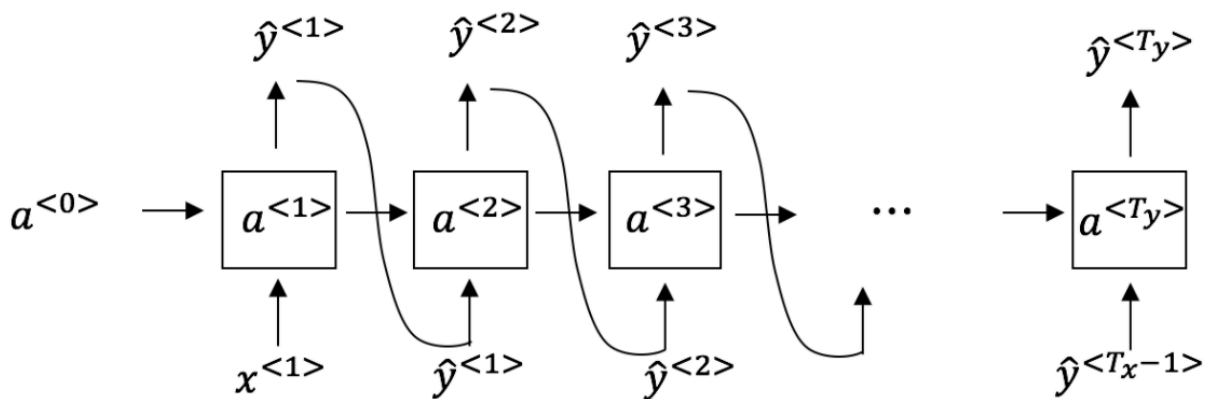
Question 5

1

point

5. Question 5

You have finished training a language model RNN and are using it to sample random sentences, as follows:



What are you doing at each time step t ?

- ☐ (i) Use the probabilities output by the RNN to pick the highest probability word for that time-step as $y^{\{<t>\}}$. (ii) Then pass the ground-truth word from the training set to the next time-step.
- ☐ (i) Use the probabilities output by the RNN to randomly sample a chosen word for that time-step as $y^{\{<t>\}}$. (ii) Then pass the ground-truth word from the training set to the next time-step.
- ☐ (i) Use the probabilities output by the RNN to pick the highest probability word for that time-step as $y^{\{<t>\}}$. (ii) Then pass this selected word to the next time-step.
- ☒ (i) Use the probabilities output by the RNN to randomly sample a chosen word for that time-step as $y^{\{<t>\}}$. (ii) Then pass this selected word to the next time-step.

Question 6

1

+

point

6. Question 6

You are training an RNN, and find that your weights and activations are all taking on the value of NaN ("Not a Number"). Which of these is the most likely cause of this problem?

☐

Vanishing gradient problem.

☒

Exploding gradient problem.

☐

ReLU activation function $g(\cdot)$ used to compute $g(z)$, where z is too large.

☐

Sigmoid activation function $g(\cdot)$ used to compute $g(z)$, where z is too large.

Question 7

1

point

7. Question 7

Suppose you are training a LSTM. You have a 10000 word vocabulary, and are using an LSTM with 100-dimensional activations $a^{<t>}$. What is the dimension of Γ_u at each time step?

☐

1

☒

100

☐

300

☐

10000

Question 8

1

point

8. Question 8

Here're the update equations for the GRU.

GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$

Alice proposes to simplify the GRU by always removing the Γ_u . I.e., setting $\Gamma_u = 1$. Betty

proposes to simplify the GRU by removing the Γ_r . I. e., setting $\Gamma_r = 1$ always. Which of these models is more likely to work without vanishing gradient problems even when trained on very long input sequences?

- ☐ Alice's model (removing Γ_u), because if $\Gamma_r \approx 0$ for a timestep, the gradient can propagate back through that timestep without much decay.
- ☐ Alice's model (removing Γ_u), because if $\Gamma_r \approx 1$ for a timestep, the gradient can propagate back through that timestep without much decay.
- ☒ Betty's model (removing Γ_r), because if $\Gamma_u \approx 0$ for a timestep, the gradient can propagate back through that timestep without much decay.
- ☐ Betty's model (removing Γ_r), because if $\Gamma_u \approx 1$ for a timestep, the gradient can propagate back through that timestep without much decay.

Question 9

1

point

9. Question 9

Here are the equations for the GRU and the LSTM:

GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$

LSTM

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * c^{<t>}$$

From these, we can see that the Update Gate and Forget Gate in the LSTM play a role similar to _____ and _____ in the GRU. What should go in the the blanks?

- ☒ Γ_u and $1 - \Gamma_u$
- ☐ Γ_u and Γ_r
- ☐ $1 - \Gamma_u$ and Γ_u
- ☐ Γ_r and Γ_u

Question 10

1

point

10. Question 10

You have a pet dog whose mood is heavily dependent on the current and past few days' weather. You've collected data for the past 365 days on the weather, which you represent as a sequence as $x_{<1>}, \dots, x_{<365>}$. You've also collected data on your dog's mood, which you represent as $y_{<1>}, \dots, y_{<365>}$. You'd like to build a model to map from $x \rightarrow y$. Should you use a Unidirectional RNN or Bidirectional RNN for this problem?

- ☐ Bidirectional RNN, because this allows the prediction of mood on day t to take into account more information.
- ☐ Bidirectional RNN, because this allows backpropagation to compute more accurate gradients.
- ☒ Unidirectional RNN, because the value of $y_{<t>}$ depends only on $x_{<1>}, \dots, x_{<t>}$, but not on $x_{<t+1>}, \dots, x_{<365>}$
- ☐ Unidirectional RNN, because the value of $y_{<t>}$ depends only on $x_{<t>}$, and not other days' weather.

1. Question 1

Suppose you learn a word embedding for a vocabulary of 10000 words. Then the embedding vectors should be 10000 dimensional, so as to capture the full range of variation and meaning in those words.

- ☐ True
- ☒ False

Question 2

1

point

2. Question 2

What is t-SNE?

- ☐ A linear transformation that allows us to solve analogies on word vectors
- ☒ A non-linear dimensionality reduction technique
- ☐ A supervised learning algorithm for learning word embeddings
- ☐ An open-source sequence modeling library

Question 3

1

point

3. Question 3

Suppose you download a pre-trained word embedding which has been trained on a huge corpus of text. You then use this word embedding to train an RNN for a language task of recognizing if someone is happy from a short snippet of text, using a small training set.

x (input text)	y (happy?)
I'm feeling wonderful today!	1
I'm bummed my cat is ill.	0
Really enjoying this!	1

Then even if the word “ecstatic” does not appear in your small training set, your RNN might reasonably be expected to recognize “I’m ecstatic” as deserving a label $y = 1$.



True

False

Question 4

1

point

4. Question 4

Which of these equations do you think should hold for a good word embedding? (Check all that apply)



$e_{\text{boy}} - e_{\text{girl}} \approx e_{\text{brother}} - e_{\text{sister}}$ *eboy-egirl≈ebrother-esister*

$e_{\text{boy}} - e_{\text{girl}} \approx e_{\text{sister}} - e_{\text{brother}}$ *eboy-egirl≈esister-ebrother*

$e_{\text{boy}} - e_{\text{brother}} \approx e_{\text{girl}} - e_{\text{sister}}$ *eboy-ebrother≈egirl-esister*

$e_{\text{boy}} - e_{\text{brother}} \approx e_{\text{sister}} - e_{\text{girl}}$ *eboy-ebrother≈esister-egirl*

Question 5

1

point

5. Question 5

Let E be an embedding matrix, and let o_{1234} be a one-hot vector corresponding to word 1234. Then to get the embedding of word 1234, why don't we call $E * o_{1234}$ $E*o_{1234}$ in Python?



It is computationally wasteful.

The correct formula is $EAT * o_{1234}$ $E^T * o_{1234}$

- ☐ The correct formula is $E_{t \sim P(\{1234\})} E_t * O_{1234}$.
- ☐ This doesn't handle unknown words (<UNK>).
- ☐ None of the above: calling the Python snippet as described above is fine.

Question 6

1

point

6. Question 6

When learning word embeddings, we create an artificial task of estimating $P(\text{target}|\text{context})$. It is okay if we do poorly on this artificial prediction task; the more important by-product of this task is that we learn a useful set of word embeddings.

- ☒ True
- ☐ False

Question 7

1

point

7. Question 7

In the word2vec algorithm, you estimate $P(t|c)$, where tt is the target word and cc is a context word. How are tt and cc chosen from the training set? Pick the best answer.

- ☐ cc is the one word that comes immediately before tt .
- ☐ cc is a sequence of several words immediately before tt .
- ☒ cc and tt are chosen to be nearby words.
- ☐ cc is the sequence of all the words in the sentence before tt .

Question 8

1

point

8. Question 8

Suppose you have a 10000 word vocabulary, and are learning 500-dimensional word embeddings. The word2vec model uses the following softmax function:

$$P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{t=1}^{10000} e^{\theta_t^T e_c}}$$

Which of these statements are correct? Check all that apply.

- ☒ θ_t and e_c are both 500 dimensional vectors.
- ☐ θ_t and e_c are both 10000 dimensional vectors.
- ☐ θ_t and e_c are both trained with an optimization algorithm such as Adam or

- ☒ $\theta_{t\theta}$ and e_{cc} are both trained with an optimization algorithm such as Adam or gradient descent.
- ☐ After training, we should expect $\theta_{t\theta}$ to be very close to e_{cc} when tt and cc are the same word.

Question 9

1

point

9. Question 9

Suppose you have a 10000 word vocabulary, and are learning 500-dimensional word embeddings. The GloVe model minimizes this objective:

$$\min \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij})(\theta_i e_j + b_i + b_j - \log X_{ij})^2$$

Which of these statements are correct? Check all that apply.

- ☐ θ_i and e_j should be initialized to 0 at the beginning of training.
- ☒ θ_i and e_j should be initialized randomly at the beginning of training.
- ☒ X_{ij} is the number of times word i appears in the context of word j .
- ☒ The weighting function $f(\cdot)$ must satisfy $f(0) = 0$.

Question 10

1

point

10. Question 10

You have trained word embeddings using a text dataset of m_1 words. You are considering using these word embeddings for a language task, for which you have a separate labeled dataset of m_2 words. Keeping in mind that using word embeddings is a form of transfer learning, under which of these circumstance would you expect the word embeddings to be helpful?

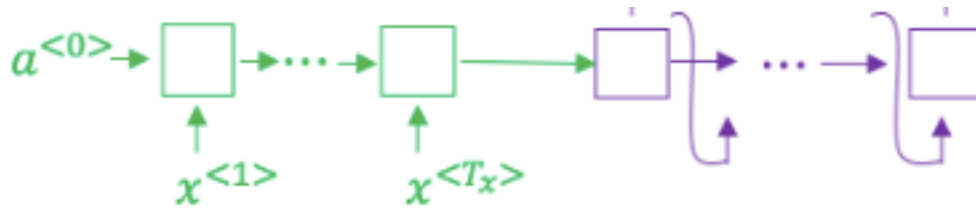
- ☒ $m_1 \gg m_2$
- ☐ $m_1 \ll m_2$

1. Question 1

Consider using this encoder-decoder model for machine translation.

$\hat{y}^{<1>}$

$\hat{y}^{<T_y>}$



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 xx .

☐
☒

True

False

Question 2

1

point

2. Question 2

In beam search, if you increase the beam width BB , which of the following would you expect to be true? Check all that apply.

☒
☒

Beam search will run more slowly.

Beam search will use up more memory.

☒

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

☐

Beam search will converge after fewer steps.

Question 3

1

point

3. Question 3

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

☒
☐

True

False

Question 4

1

point

4. Question 4

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

On a dev set example, given an input audio clip, your algorithm outputs the

transcript y^{\wedge} = "I'm building an A Eye system in Silly con Valley.", whereas a human gives a much superior transcript $y^{\wedge*} = y_*$ = "I'm building an AI system in Silicon Valley."

According to your model,

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

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

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

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

Question 5

1

point

5. Question 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, $P(y_*|x) > P(y^{\wedge}|x)$. This suggest you should focus your attention on improving the search algorithm.

- ☒ True.
- ☐ False.

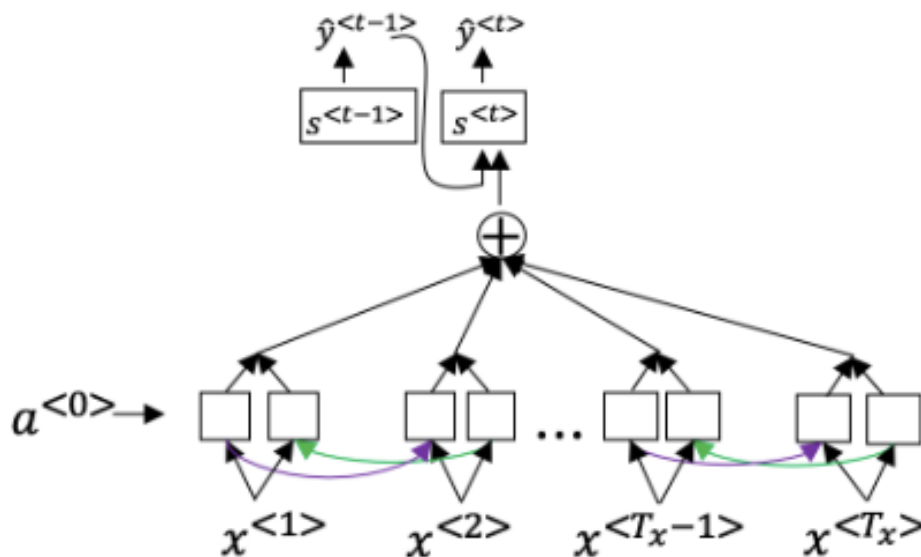
Question 6

1

point

6. Question 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.

- ☒ 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.)
- ☐ 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.)
- ☐ $\sum_t \alpha_{<t,t'>} = 1$ (Note the summation is over t .)
- ☒ $\sum_{t'} \alpha_{<t,t'>} = 1$ (Note the summation is over t' .)

Question 7

1

point

7. Question 7

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

computed using a small neural network:

We can't replace $s^{<t-1>}s_{<t-1>}$ with $s^{<t>}s_{<t>}$ as an input to this neural network. This is because $s^{<t>}s_{<t>}$ 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 $s^{<t>}s_{<t>}$ yet.



True



False

Question 8

1

point

8. Question 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:



The input sequence length $T_x T_x$ is large.



The input sequence length $T_x T_x$ is small.

Question 9

1

point

9. Question 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__o__o__kk__b__o__o__o__o__o__o__o__o__o__kkk`



cokbok



cookbook



cook book



coookkboooooookkk

Question 10

1

point

10. Question 10

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



Features of the audio (such as spectrogram features) at time tt .



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



Whether the trigger word is being said at time tt .



Whether someone has just finished saying the trigger word at time tt .