Congratulations! You passed! TO PASS 80% or higher

Recurrent Neural Networks

LATEST SUBMISSION GRADE 100%

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

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?

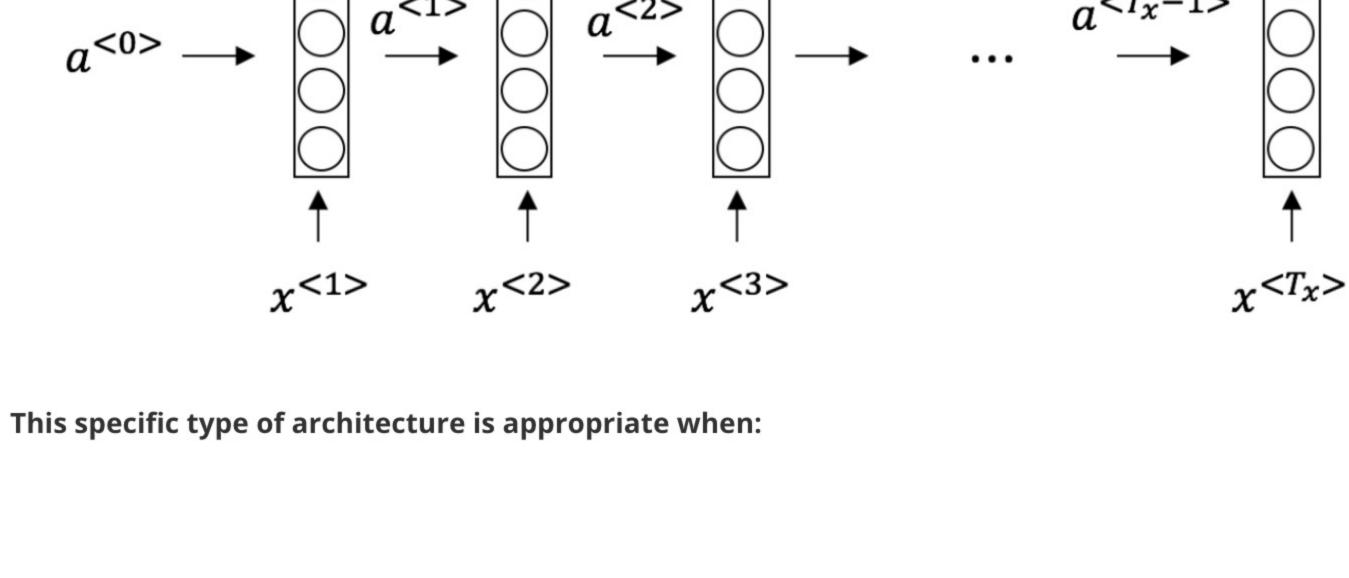
 $x^{< i>(j)}$

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

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

We index into the i^{th} row first to get the i^{th} training example (represented by parentheses), then the j^{th} column to get the j^{th} word (represented by the brackets).

Consider this RNN:



 $\bigcap T_x = 1$

It is appropriate when every input should be matched to an output.

 $\bigcap T_x > T_y$

Correct

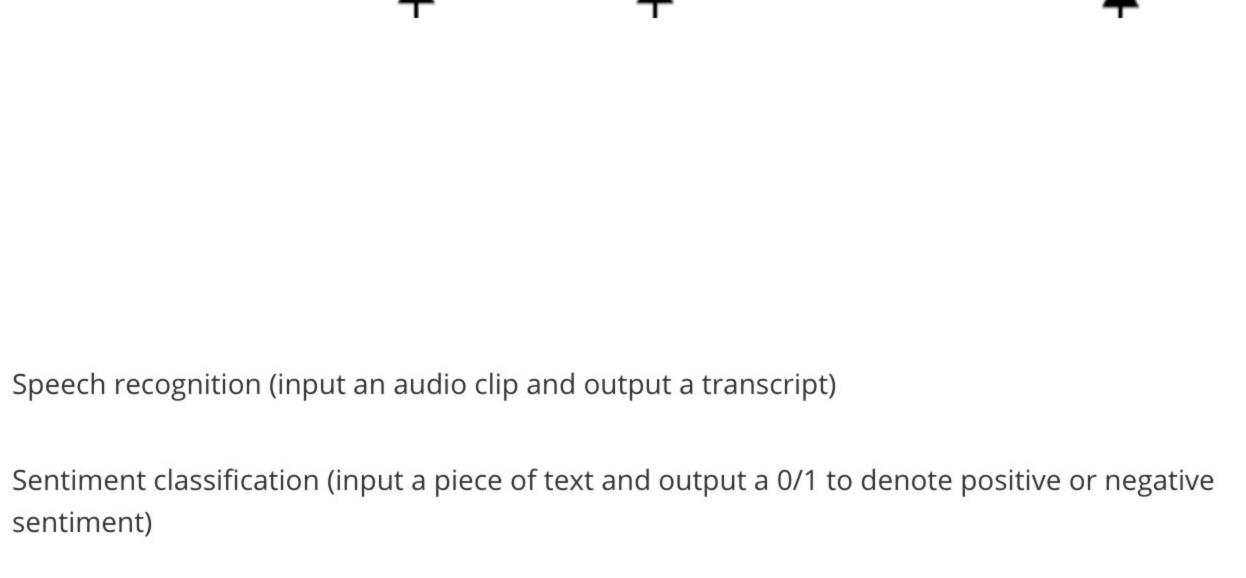
Correct

gender)

Correct

Correct!

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



Gender recognition from speech (input an audio clip and output a label indicating the speaker's

You are training this RNN language model.

Image classification (input an image and output a label)

At the t^{th} time step, what is the RNN doing? Choose the best answer. Estimating $P(y^{<1>}, y^{<2>}, \dots, y^{<t-1>})$

Yes, in a language model we try to predict the next step based on the knowledge of all prior

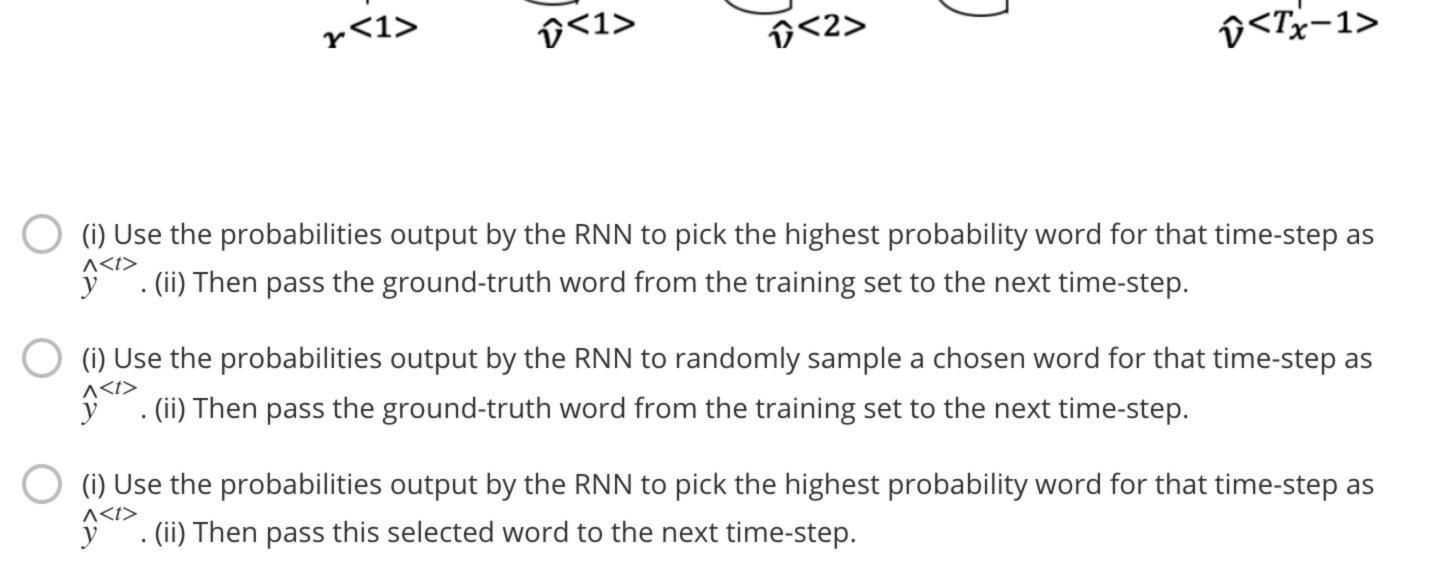
5. You have finished training a language model RNN and are using it to sample random sentences, as

✓ Correct

follows:

steps.

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



(i) Use the probabilities output by the RNN to randomly sample a chosen word for that time-step as

Correct Yes!

 $\hat{y}^{< t>}$. (ii) Then pass this selected word to the next time-step.

Vanishing gradient problem.

Correct

- 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?
- Sigmoid activation function g(.) used to compute g(z), where z is too large.

Suppose you are training a LSTM. You have a 10000 word vocabulary, and are using an LSTM with

Correct, Γ_u is a vector of dimension equal to the number of hidden units in the LSTM.

100 300

100-dimensional activations $a^{< t>}$. What is the dimension of Γ_u at each time step?

✓ Correct

10000

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

GRU

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

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

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

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

to work without vanishing gradient problems even when trained on very long input sequences?

simplify the GRU by removing the Γ_r . I. e., setting Γ_r = 1 always. Which of these models is more likely

Alice's model (removing Γ_u), because if $\Gamma_r pprox 0$ for a timestep, the gradient can propagate back through that timestep without much decay. Alice's model (removing Γ) hereuse if $\Gamma\sim 1$ for a timesten the gradient can propagate back

through that timestep without much decay.

through that timestep without much decay.

Here are the equations for the GRU and the LSTM:

Correct Yes. For the signal to backpropagate without vanishing, we need $c^{< t>}$ to be highly dependant on $c^{< t-1>}$.

Betty's model (removing Γ_r), because if $\Gamma_u \approx 1$ for a timestep, the gradient can propagate back

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

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?

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

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

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

Correct Yes, correct!

 \bigcap Γ_r and Γ_u

 \bigcirc Γ_u and $1 - \Gamma_u$

 $y^{<1>},\ldots,y^{<365>}$. You'd like to build a model to map from x o 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.

 $x^{<1>}, \dots, x^{<365>}$. You've also collected data on your dog's mood, which you represent as

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

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.

Keep Learning

GRADE

100%

1/1 point

1/1 point

1 / 1 point

1/1 point

1 / 1 point

1/1 point

1/1 point

1/1 point

1/1 point