✓ **Congratulations! You passed!**

| Grade | Latest Submission | To pass 80% or | |
|---|---|---|---|
| received 100% | Grade 100% | higher | **Go to next item** |

**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?

1 / 1 point

- ⦿ $x^{(i)<j>}$
- ○ $x^{<i>(j)}$
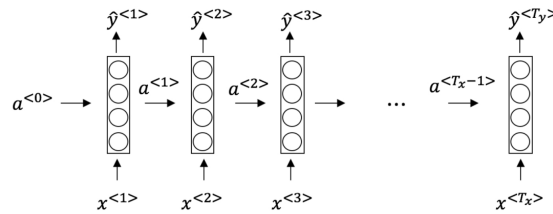- ○ $x^{(j)<i>}$
- ○ $x^{<j>(i)}$

⟋ **Expand**

⊘ **Correct**
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).

**2.** Consider this RNN:

1 / 1 point



True/False: This specific type of architecture is appropriate when Tx=Ty
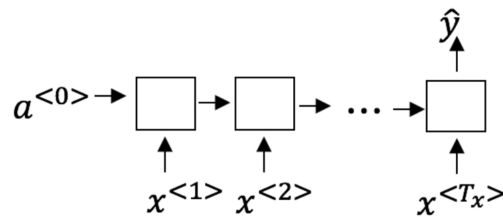
- ⦿ True
- ○ False

⟋ **Expand**

⊘ **Correct**
It is appropriate when the input sequence and the output sequence have the same length or size.

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

1 / 1 point



- ☐ 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)

  ✓ **Correct**
  Correct!

- ☐ 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)
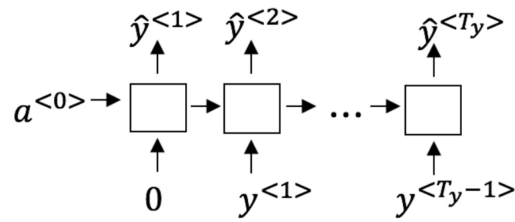
  ✓ **Correct**
  Correct!

⟋ **Expand**

⊘ **Correct**
Great, you got all the right answers.

**4.** You are training this RNN language model.



At the $t^{th}$ time step, what is the RNN doing?

○ Estimating $P(y^{<1>}, y^{<2>}, \ldots, y^{<t-1>})$

○ Estimating $P(y^{<t>})$

Typesetting math: 100% $^{<t>} \mid y^{<1>}, y^{<2>}, \ldots, y^{<t-1>})$
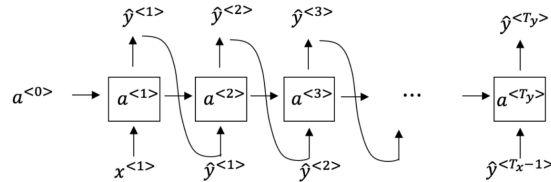
⬈ **Expand**

✓ **Correct**
Yes, in a language model we try to predict the next step based on the knowledge of all prior steps.

---

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

لقد انتهيت من تدريب نموذج اللغة RNN وتستخدمه لأخذ عينات من الجمل العشوائية ، على النحو التالي:

True/False: In this sample sentence, step t uses the probabilities output by the RNN to randomly sample a chosen word for that time-step. Then it passes this selected word to the next time-step.

صواب/خطأ: في هذه الجملة النموذجية، تستخدم الخطوة t بواسطة RNN الاحتمالات الناتجة لتلك الخطوة الزمنية. ثم لأخذ عينة عشوائية من كلمة مختارة لتلك الخطوة الزمنية. يمرر هذه الكلمة المحددة إلى الخطوة الزمنية التالية.

○ False

◉ True

⬈ **Expand**

✓ **Correct**
Step t uses the probabilities output by the RNN to randomly sample a chosen word for that time-step. Then it passes this selected word to the next time-step.

---

**6.** True/False: If you are training an RNN model, and find that your weights and activations are all taking on the value of NaN ("Not a Number") then you have an exploding gradient problem.

○ False

◉ True

⬈ **Expand**

✓ **Correct**
Correct! Exploding gradients happen when large error gradients accumulate and result in very large updates to the NN model weights during training. These weights can become too large and cause an overflow, identified as NaN.

---

**7.** Suppose you are training an LSTM. You have a 50000 word vocabulary, and are using an LSTM with 500-dimensional activations $a^{<t>}$.

مع عمليات تنشيط 500 الأبعاد LSTM لديك مفردات 50000 كلمة ، وتستخدم .LSTM لنفترض أنك تقوم بتدريب a. What is the dimension of $\Gamma_u$ at each time step?

◉ 500

○ 200

○ 50000

○ 5

<button>Expand</button>

✓ **Correct**
Correct, $\Gamma_u$ is a vector of dimension equal to the number of hidden units in the LSTM.

---

**8.** Here are 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 $\Gamma_u$. I.e., setting $\Gamma_u$ = 0. Betty proposes to simplify the GRU by removing the $\Gamma_r$. 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 $\Gamma_u$), because if $\Gamma_r \approx 0$ for a timestep, the gradient can propagate back through that timestep without much decay.

○ Alice's model (removing $\Gamma_u$), because if $\Gamma_r \approx 1$ for a timestep, the gradient can propagate back through that timestep without much decay.

◉ Betty's model (removing $\Gamma_r$), because if $\Gamma_u \approx 0$ for a timestep, the gradient can propagate back through that timestep without much decay.

○ Betty's model (removing $\Gamma_r$), because if $\Gamma_u \approx 1$ for a timestep, the gradient can propagate back through that timestep without much decay.

<button>Expand</button>

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

**1 / 1 point**

---

**9.** Here are the equations for the GRU and the LSTM:

| GRU | LSTM |
|---|---|
| $$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$ | $$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$ |
| $$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$ | $$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$ |
| $$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$ | $$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$ |
| $$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$ | $$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$ |
| $$a^{<t>} = c^{<t>}$$ | $$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$ |
| | $$a^{<t>} = \Gamma_o * \tanh c^{<t>}$$ |

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 blanks?

◉ $\Gamma_u$ and $1 - \Gamma_u$

○ $\Gamma_u$ and $\Gamma_r$

○ $1 - \Gamma_u$ and $\Gamma_u$

○ $\Gamma_r$ and $\Gamma_u$

<button>Expand</button>

✓ **Correct**
Yes, correct!

**1 / 1 point**

---

**10.** Your 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 $x^{<1>}, \ldots, x^{<365>}$

<span dir="rtl">You've .x يعتمد مزاجك بشكل كبير على الطقس الحالي والأيام القليلة الماضية. لقد جمعت بيانات عن آخر 365 يوما عن الطقس، والتي تمثلها كتسلسل ك</span> also collected data on your mood, which you represent as $y^{<1>}, \ldots, y^{<365>}$. You'd like to build a model to map from x→y. Should you use a Unidirectional RNN or Bidirectional RNN for this problem?

◉ Unidirectional RNN, because the value of $y^{<t>}$ depends only on $x^{<1>}, \ldots, x^{<t>}$, but not on $x^{<1>}, \ldots, x^{<365>}$.

○ Bidirectional RNN, because this allows the prediction of mood on day t to take into account more information.

○ Unidirectional RNN, because the value of $y^{<t>}$ depends only on $x^{<t>}$, and not other days' weather.

○ Bidirectional RNN, because this allows backpropagation to compute more accurate gradients.

**1 / 1 point**

Expand

Correct

Expand

Correct