## **~**

# **Congratulations! You passed!**

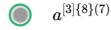
Next Item



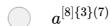
Which notation would you use to denote the 3rd layer's activations when the input is the 7th example from the 8th minibatch?

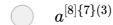
1/1 point

 $a^{[3]\{7\}(8)}$ 



Correct

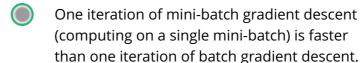




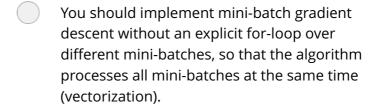


2. Which of these statements about mini-batch gradient descent do you agree with?

1/1 point



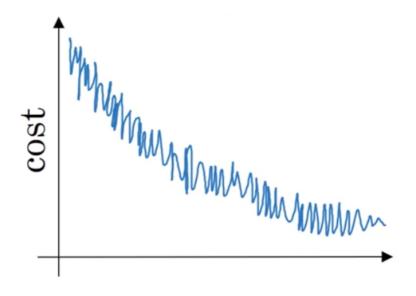




		using batch gradient descent.		
<b>~</b>	3.	Why is the best mini-batch size usually not 1 and not m, but instead something in-between?		
1/1 point		If the mini-batch size is 1, you end up having to process the entire training set before making any progress.		
		Un-selected is correct		
		If the mini-batch size is 1, you lose the benefits of vectorization across examples in the mini-batch.		
		Correct		
		If the mini-batch size is m, you end up with stochastic gradient descent, which is usually slower than mini-batch gradient descent.		
	Un-selected is correct			
		If the mini-batch size is m, you end up with batch gradient descent, which has to process the whole training set before making progress.		
		Correct		

Training one epoch (one pass through the training set) using mini-batch gradient descent is faster than training one epoch

Suppose your learning algorithm's cost J, plotted as a function of the number of iterations, looks like this:



Which of the following do you agree with?

- Whether you're using batch gradient descent or mini-batch gradient descent, something is wrong.
- Whether you're using batch gradient descent or mini-batch gradient descent, this looks acceptable.
- If you're using mini-batch gradient descent, this looks acceptable. But if you're using batch gradient descent, something is wrong.

### Correct

If you're using mini-batch gradient descent, something is wrong. But if you're using batch gradient descent, this looks acceptable.



5. Suppose the temperature in Casablanca over the first three days of January are the same:

1/1 point

Jan 1st:  $heta_1=10^oC$ 

Jan 2nd:  $heta_2 10^o C$ 

(We used Fahrenheit in lecture, so will use Celsius here in honor of the metric world.)

Say you use an exponentially weighted average with  $\beta=0.5$  to track the temperature:  $v_0=0$ ,  $v_t=\beta v_{t-1}+(1-\beta)\theta_t$ . If  $v_2$  is the value computed after day 2 without bias correction, and  $v_2^{corrected}$  is the value you compute with bias correction. What are these values? (You might be able to do this without a calculator, but you don't actually need one. Remember what is bias correction doing.)

$$v_2=10$$
,  $v_2^{corrected}=10$ 

$$v_2=10$$
,  $v_2^{corrected}=7.5$ 

$$v_2=7.5$$
,  $v_2^{corrected}=7.5$ 

$$v_2=7.5$$
,  $v_2^{corrected}=10$ 

Correct



6. Which of these is NOT a good learning rate decay scheme? Here, t is the epoch number.

1/1 point

$$lpha = rac{1}{1+2*t}lpha_0$$

$$lpha = rac{1}{\sqrt{t}} lpha_0$$

$$lpha = 0.95^tlpha_0$$

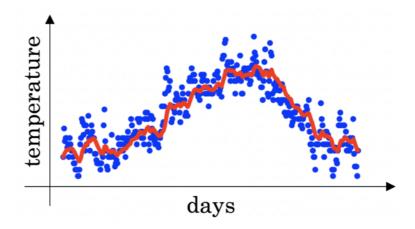
$$\bigcirc \quad \alpha = e^t \alpha_0$$

Correct



7. You use an exponentially weighted average on the London temperature dataset. You use the following to track the temperature:  $v_t = \beta v_{t-1} + (1-\beta)\theta_t$ . The

red line below was computed using  $\beta=0.9$ . What would happen to your red curve as you vary  $\beta$ ? (Check the two that apply)



Decreasing eta will shift the red line slightly to the right.

#### **Un-selected is correct**

Increasing  $\beta$  will shift the red line slightly to the right.

#### Correct

True, remember that the red line corresponds to  $\beta=0.9$ . In lecture we had a green line \$\$\beta=0.98\$) that is slightly shifted to the right.

Decreasing  $\beta$  will create more oscillation within the red line.

#### Correct

True, remember that the red line corresponds to  $\beta=0.9$ . In lecture we had a yellow line \$\$\beta=0.98\$ that had a lot of oscillations.

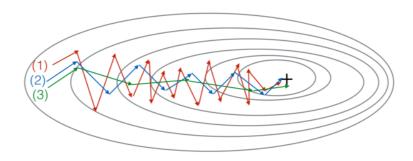
Increasing $eta$ will create more oscillations
within the red line.

#### **Un-selected is correct**



8. Consider this figure:

1/1 point



These plots were generated with gradient descent; with gradient descent with momentum ( $\beta$  = 0.5) and gradient descent with momentum ( $\beta$  = 0.9). Which curve corresponds to which algorithm?



(1) is gradient descent. (2) is gradient descent with momentum (small  $\beta$ ). (3) is gradient descent with momentum (large  $\beta$ )



- (1) is gradient descent. (2) is gradient descent with momentum (large  $\beta$ ) . (3) is gradient descent with momentum (small  $\beta$ )
- (1) is gradient descent with momentum (small  $\beta$ ), (2) is gradient descent with momentum (small  $\beta$ ), (3) is gradient descent
- (1) is gradient descent with momentum (small  $\beta$ ). (2) is gradient descent. (3) is gradient descent with momentum (large  $\beta$ )

1/1 point

Suppose batch gradient descent in a deep network is taking excessively long to find a value of the parameters that achieves a small value for the cost function  $\mathcal{J}(W^{[1]},b^{[1]},...,W^{[L]},b^{[L]})$ . Which of the following techniques could help find parameter values that attain a small value for  $\mathcal{J}$ ? (Check all that apply)

- 1	-		-	n
	п			п
				ш
				ш
				ш

Try using Adam





Try better random initialization for the weights





Try mini-batch gradient descent

Correct



Try tuning the learning rate lpha





Try initializing all the weights to zero

**Un-selected is correct** 



1 Which of the following statements about Adam is False?

1/1 point

We usually use "default" values for the hyperparameters  $eta_1, eta_2$  and arepsilon in Adam (  $eta_1=0.9, eta_2=0.999, arepsilon=10^{-8}$  )

	The learning rate hyperparameter $lpha$ in Adam usually needs to be tuned.					
	Adam should be used with batch gradient computations, not with mini-batches.					
Correct						
	Adam combines the advantages of RMSProp and momentum					

