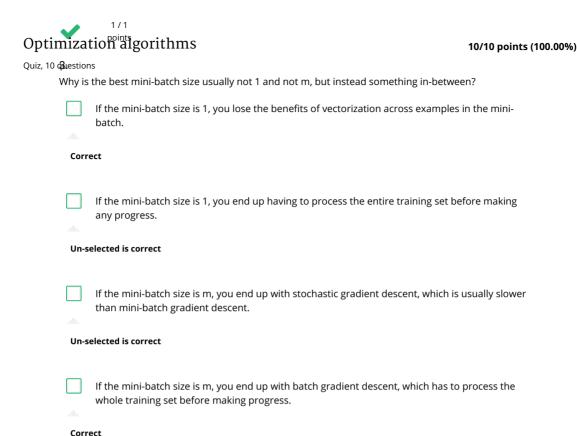
otimization algorithms , 10 questions		10/10 points (100.00%)	
✓ Co	ongratulations! You passed!	Next Item	
~	1 / 1 points		
	notation would you use to denote the 3rd layer's activations when the inpune 8th minibatch?	it is the 7th example	
\bigcirc	$a^{[8]\{7\}(3)}$		
\bigcirc	$a^{[8]\{3\}(7)}$		
\bigcirc	$a^{[3]\{7\}(8)}$		
0	$a^{[3]\{8\}(7)}$		
Corr	ect		
~	1 / 1 points		
2. Which	of these statements about mini-batch gradient descent do you agree with?		
0	One iteration of mini-batch gradient descent (computing on a single minione iteration of batch gradient descent.		
Corr	ect		
0	Training one epoch (one pass through the training set) using mini-batch g faster than training one epoch using batch gradient descent.	radient descent is	
\bigcirc	You should implement mini-batch gradient descent without an explicit for mini-batches, so that the algorithm processes all mini-batches at the same		

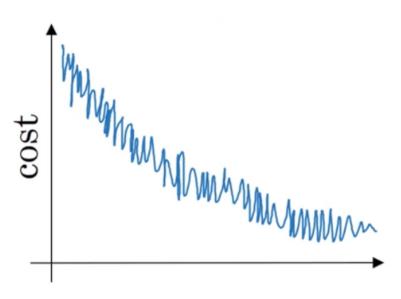
 \leftarrow



10/10 points (100.00%)

Quiz, 10 4uestions

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?

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

10/10 points (100.00%)

Quiz, 10 duestions

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

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.)

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

$$igcup_2=7.5$$
, $v_2^{corrected}=10$

Correct

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

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



1/1 points

6.

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

$$\alpha = \frac{1}{\sqrt{t}} \alpha_0$$

$$\bigcirc \quad \alpha = \frac{1}{1+2*t}\alpha_0$$

$$\alpha = 0.95^t lpha_0$$

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

Correct

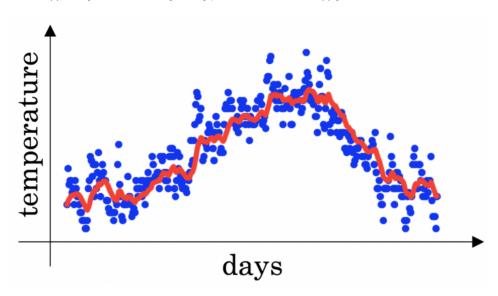
4 of 8

Un-selected is correct

10/10 points (100.00%)

Quiz, 10 duestions

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 β ? (Check the two that apply)



Un-selected is correct

Un-selected is correct

Increasing β 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 β will create more oscillation within the red line.

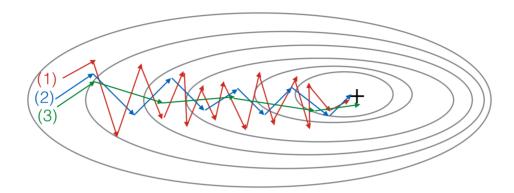
Correct

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.

10/10 points (100.00%)

Quiz, 10 **B**uestions

Consider this figure:



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

(1) is gradient descent. (2) is gradient descent with momentum (large β) . (3) is gradient descent with momentum (small β)

(1) is gradient descent. (2) is gradient descent with momentum (small β). (3) is gradient descent with momentum (large β)

Correct

(1) is gradient descent with momentum (small β), (2) is gradient descent with momentum (small β), (3) is gradient descent

(1) is gradient descent with momentum (small β). (2) is gradient descent. (3) is gradient descent with momentum (large β)



10/10 points (100.00%)

Quiz, 10 Questions

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)

following techniques could help find parameter values that attain a small value for \mathcal{J} ? (Check all that apply)
Try using Adam
Correct
Try better random initialization for the weights
Correct
Try tuning the learning rate $lpha$
Correct
Try initializing all the weights to zero
Un-selected is correct
Try mini-batch gradient descent
Correct
1/1 points
10.
Which of the following statements about Adam is False?
We usually use "default" values for the hyperparameters eta_1,eta_2 and $arepsilon$ in Adam ($eta_1=0.9$, $eta_2=0.999$, $arepsilon=10^{-8}$)
Adam should be used with batch gradient computations, not with mini-batches.
Correct
Adam combines the advantages of RMSProp and momentum

