

Congratulations! You passed!

Grade
received **100%**

Latest Submission
Grade 100%

To pass 80% or
higher

[Go to next item](#)

1. 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

- ☐ $a^{[3]}(7)(8)$
☐ $a^{[8]}(7)(3)$
☐ $a^{[8]}(3)(7)$
☒ $a^{[3]}(8)(7)$

 [Expand](#)

 **Correct**

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

1 / 1 point

- ☐ Training one epoch (one pass through the training set) using mini-batch gradient descent is faster than training one epoch using batch gradient descent.
☐ You should implement mini-batch gradient descent without an explicit for-loop over different mini-batches so that the algorithm processes all mini-batches at the same time (vectorization).
☒ When the mini-batch size is the same as the training size, mini-batch gradient descent is equivalent to batch gradient descent.

 [Expand](#)

 **Correct**

Correct. Batch gradient descent uses all the examples at each iteration, this is equivalent to having only one mini-batch of the size of the complete training set in mini-batch gradient descent.

3. Which of the following is true about batch gradient descent?

1 / 1 point

- ☒ It is the same as the mini-batch gradient descent when the mini-batch size is the same as the size of the training set.
☐ It has as many mini-batches as examples in the training set.
☐ It is the same as stochastic gradient descent, but we don't use random elements.

 [Expand](#)

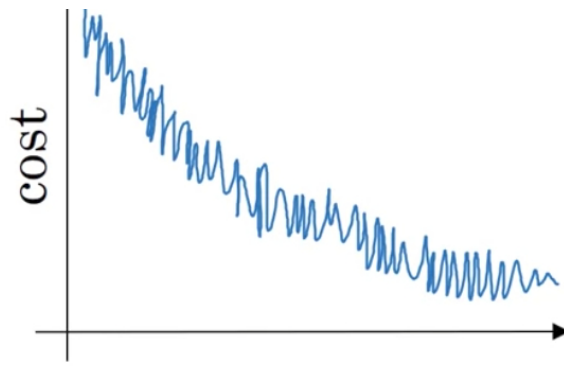
 **Correct**

Correct. When using batch gradient descent there is only one mini-batch thus it is equivalent to batch gradient descent.

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

1 / 1 point





Which of the following do you agree with?

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

[Expand](#)

✓ Correct

5. Suppose the temperature in Casablanca over the first two days of March are the following:

1 / 1 point

March 1st: $\theta_1 = 30^\circ \text{ C}$

March 2nd: $\theta_2 = 15^\circ \text{ C}$

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 $v_2^{\text{corrected}}$ is the value you compute with bias correction. What are these values?

- ☐ $v_2 = 20, v_2^{\text{corrected}} = 15$.
- ☐ $v_2 = 15$

[Expand](#)

✓ Correct

Correct. $v_2 = \beta v_{t-1} + (1 - \beta) \theta_t$ thus $v_1 = 15, v_2 = 15$. Using the bias correction $\frac{v_t}{1 - \beta^t}$ we get $\frac{15}{1 - (0.5)^2} = 20$.

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

1 / 1 point

- ☐ $\alpha = \frac{\alpha_0}{1 + 3t}$
- ☐ $\alpha = e^{-0.01t} \alpha_0$.
- ☐ $\alpha = \frac{\alpha_0}{\sqrt{1 + t}}$.
- ☒ $\alpha = 1.01^t \alpha_0$

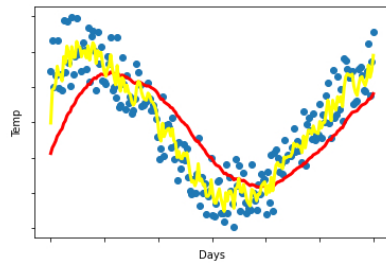
Expand

✓ Correct

Correct. This is not a good learning rate decay since it is an increasing function of t .

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 yellow and red lines were computed using values β_1 and β_2 respectively. Which of the following are true?

1 / 1 point



- ☐ $\beta_1 > \beta_2$.
- ☐ $\beta_1 = \beta_2$.
- ☒ $\beta_1 < \beta_2$.
- ☐ $\beta_1 = 0, \beta_2 > 0$.

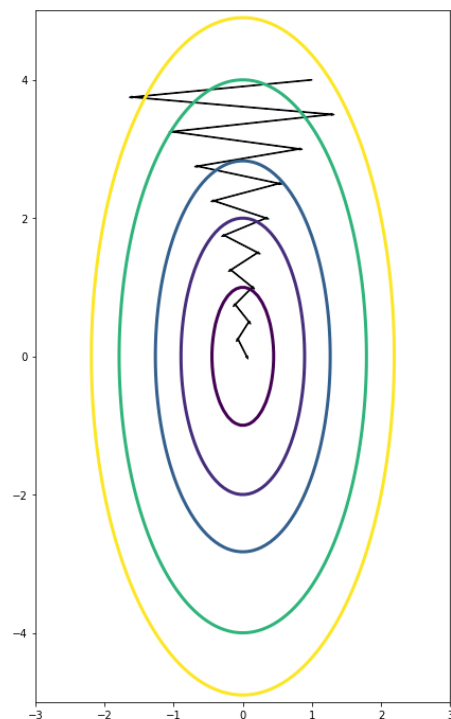
Expand

✓ Correct

Correct. $\beta_1 < \beta_2$ since the yellow curve is noisier.

8. Consider the figure:

1 / 1 point



Suppose this plot was generated with gradient descent with momentum $\beta = 0.01$. What happens if we increase the value of β to 0.1?

- ☐ The gradient descent process moves more in the horizontal and the vertical axis.
- ☐ The gradient descent process starts oscillating in the vertical direction.
- ☒ The gradient descent process moves less in the horizontal direction and more in the vertical direction.
- ☐ The gradient descent process starts moving more in the horizontal direction and less in the vertical.

 Expand

 **Correct**

Yes. The use of a greater value of β causes a more efficient process thus reducing the oscillation in the horizontal direction and moving the steps more in the vertical direction.

9. 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 $\mathcal{J}(W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]})$. Which of the following techniques could help find parameter values that attain a small value for \mathcal{J} ? (Check all that apply)

1 / 1 point

- ☒ Try tuning the learning rate α

 **Correct**

- ☒ Try using Adam

 **Correct**

- ☒ Try better random initialization for the weights

 **Correct**

- ☒ Try mini-batch gradient descent

 **Correct**

- ☐ Try initializing all the weights to zero

 Expand

 **Correct**

Great, you got all the right answers.

10. Which of the following statements about Adam is **False**?

1 / 1 point

- ☐ The learning rate hyperparameter α in Adam usually needs to be tuned.
- ☐ Adam combines the advantages of RMSProp and momentum
- ☒ Adam should be used with batch gradient computations, not with mini-batches.
- ☐ We usually use "default" values for the hyperparameters β_1, β_2 and ϵ in Adam ($\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$)

 Expand

 **Correct**