# ✓ Congratulations! You passed!

**Grade received** 80%
**Latest Submission Grade** 80%
**To pass** 80% or higher

**Go to next item**

**1.** Which notation would you use to denote the 4th layer's activations when the input is the 7th example from the 3rd mini-batch?

⤢ **Expand**

⊘ **Correct**
Yes. In general $a^{[l]\{t\}(k)}$ denotes the activation of the layer $l$ when the input is the example $k$ from the mini-batch $t$.

**1 / 1 point**

**2.** Suppose you don't face any memory-related problems. Which of the following make more use of vectorization.                                                **0 / 1 point**

◯

↗ **Expand**

⊗ **Incorrect**
No: If no memory problem is faced, batch gradient descent processes all of the training set in one pass, maximizing the use of vectorization.

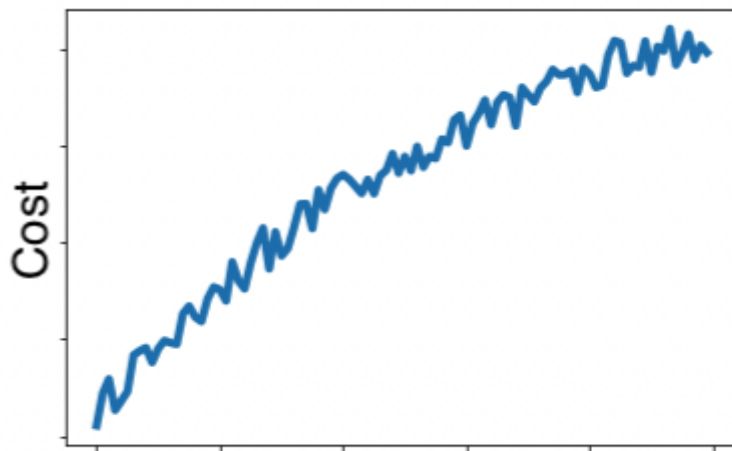**3.** Which of the following is true about batch gradient descent?                                                        **1 / 1 point**

↗ **Expand**

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

**4.** While using mini-batch gradient descent with a batch size larger than 1 but less than m     **1 / 1 point**
the plot of the cost function $J$ looks like this:



Which of the following do you agree with?

⟋  **Expand**

⊘ **Correct**
   Yes. The cost is larger than when the process started, this is not right at all.

**5.** Suppose the temperature in Casablanca over the first two days of March are the following:                          **1 / 1 point**

March 1st: $\theta_1 = 10° \text{ C}$

March 2nd: $\theta_2 = 25° \text{ C}$

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?

↗ **Expand**

⊘ **Correct**
Correct. $v_2 = \beta v_{t-1} + (1 - \beta)\,\theta_t$ thus $v_1 = 5, v_2 = 15$. Using the bias correction $\frac{v_t}{1-\beta^t}$ we get $\frac{15}{1-(0.5)^2} = 20$.

**6.** Which of the following is true about learning rate decay?                    **1 / 1 point**
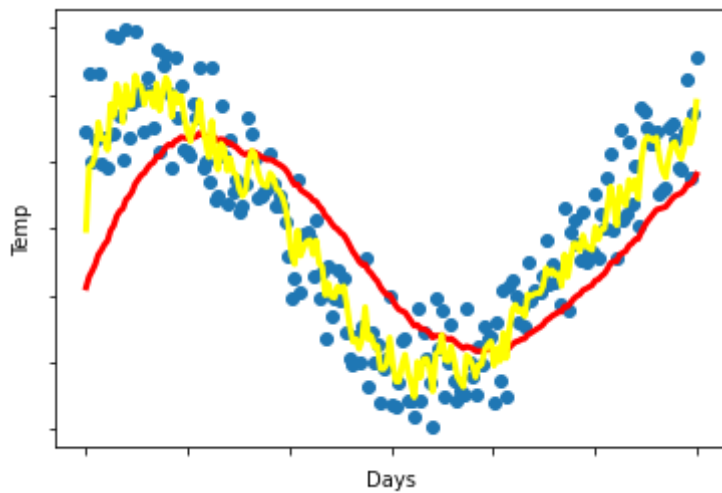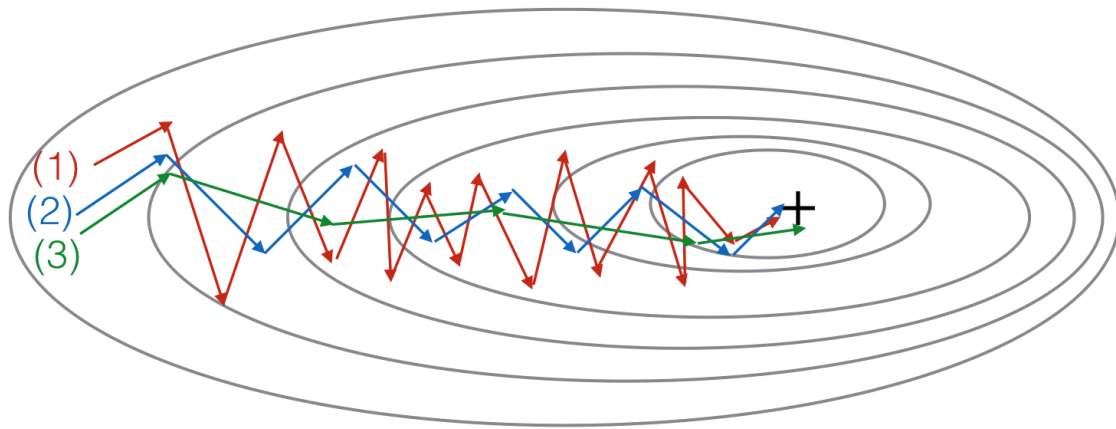
⤢ **Expand**

⊘ **Correct**

Correct. Reducing the learning rate with time reduces the oscillation around a minimum.

**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 $beta_1$ and $beta_2$ respectively. Which of the following are true?

**1 / 1 point**



↗ **Expand**

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

**8.** Consider this figure:

**0 / 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?

↗ **Expand**
↙

⊗ **Incorrect**

**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
$\mathcal{J}(W^{[1]}, b^{[1]}, ..., W^{[L]}, b^{[L]})$. 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**

↗ **Expand**

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

**10.** Which of the following statements about Adam is *False*?                    **1 / 1 point**

⤢ **Expand**

⊘ **Correct**