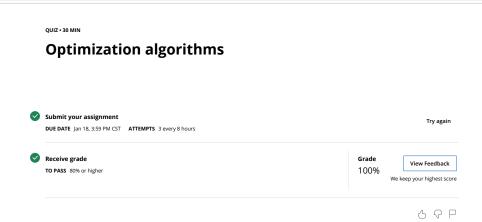
← Optimization algorithms Graded Quiz • 30 min **Due** Jan 18, 3:59 PM CST

100%



Optimization algorithms

Video: Mini-batch gradient descent
11 min

Video: Exponentially weighted averages 5 min

Video: Understanding exponentially weighted averages 9 min

Video: Gradient descent with momentum 9 min

Reading: Clarification about Upcoming Adam Optimization Video
 1 min

Video: Adam optimization algorithm
7 min

Reading: Clarification about Learning Rate Decay Video 1 min

Video: Learning rate decay 6 min

Video: The problem of local optima
5 min **Practice Questions** Quiz: Optimization algorithms
10 questions

Heroes of Deep Learning (Optional)

Video: RMSprop 7 min

Keep Learning GRADE 100% ✓ Congratulations! You passed! TO PASS 80% or higher **Optimization algorithms** LATEST SUBMISSION GRADE Which notation would you use to denote the 3rd layer's activations when the input is the 7th example from the 8th minibatch? $\bigcirc \ a^{[8]\{7\}(3)}$ $\bigcirc \ a^{[3]\{7\}\{8)}$ $\bigcirc \ a^{[8]\{3\}(7)}$ ✓ Correct 2. Which of these statements about mini-batch gradient descent do you agree with? One iteration of mini-batch gradient descent (computing on a single mini-batch) is faster than one iteration of 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). 3. Why is the best mini-batch size usually not 1 and not m, but instead something in-between? If the mini-batch size is m, you end up with stochastic gradient descent, which is usually slower than mini-batch gradient descent. If the mini-batch size is 1, you end up having to process the entire training set before making any progress. 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 If the mini-batch size is 1, you lose the benefits of vectorization across examples in the mini-batch. ✓ Correct 4. Suppose your learning algorithm's cost J, plotted as a function of the number of iterations, looks like this: 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. 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. 5. Suppose the temperature in Casablanca over the first three days of January are the same: 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_v is the value computed after day 2 without bias correction, and $v_v^{corrected}$ 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.) $igotimes v_2 = 7.5$, $v_2^{corrected} = 10$ $\bigcirc \ v_2=10, v_2^{corrected}=10$ $\bigcirc \ v_2=10$, $v_2^{corrected}=7.5$ \bigcirc $v_2=7.5$, $v_2^{corrected}=7.5$ ✓ Correct 6. Which of these is NOT a good learning rate decay scheme? Here, t is the epoch number. \bigcirc $\alpha = e^t \alpha_0$ $\alpha = \frac{1}{1+2*t}\alpha_0$ $\bigcirc \ \ lpha = 0.95^tlpha_0$ 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 β ? (Check the two that apply) temperature days igspace Increasing eta will shift the red line slightly to the right. \checkmark **Correct**True, remember that the red line corresponds to $\beta=0.9$. In lecture we had a green line \$\$\delta=0.98\$) that is slightly shifted to the right. igspace Decreasing eta will create more oscillation within the red line. 1 / 1 point

 $\hfill \Box$ Decreasing β will shift the red line slightly to the right.

 \checkmark **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.

 $\hfill \square$ Increasing β will create more oscillations within the red line.

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?

 \bigcirc (1) is gradient descent with momentum (small β). (2) is gradient descent. (3) is gradient descent with momentum (large β) igoplus (1) is gradient descent. (2) is gradient descent with momentum (small eta). (3) is gradient descent with momentum (large eta) \bigcirc (1) is gradient descent. (2) is gradient descent with momentum (large β) . (3) is gradient descent with momentum (small β)

 \bigcirc (1) is gradient descent with momentum (small eta), (2) is gradient descent with momentum (small eta), (3) is gradient descent ✓ Correct

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 $\mathcal{J}(\mathcal{C}(\mathsf{check}))$ all that apply)

Try initializing all the weights to zero Try better random initialization for the weights

✓ Correct Try mini-batch gradient descent

✓ Correct

Try using Adam

✓ Correct

✓ Correct lacksquare Try tuning the learning rate lpha

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

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

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

✓ Correct