Optimization algorithms Graded Quiz • 30 min **GRADE** Congratulations! You passed! **Keep Learning** 100% TO PASS 80% or higher **Optimization algorithms** LATEST SUBMISSION GRADE 100% Which notation would you use to denote the 3rd layer's activations when the input is the 7th 1/1 point example from the 8th minibatch? $a^{[8]\{7\}(3)}$ $a^{[8]\{3\}(7)}$ $a^{[3]\{8\}(7)}$ $a^{[3]\{7\}(8)}$ Which of these statements about mini-batch gradient descent do you agree with? 1/1 point 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 minibatches, so that the algorithm processes all mini-batches at the same time (vectorization). Training one epoch (one pass through the training set) using mini-batch gradient descent is faster than training one epoch using batch gradient descent. Correct 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 mini-batch gradient descent. 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 Suppose your learning algorithm's cost J, plotted as a function of the number of iterations, looks 1/1 point like this: Which of the following do you agree with? 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. ✓ Correct 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 eta=0.5 to track the temperature: $v_0=0$, $v_t = eta v_{t-1} + (1-eta) heta_t$. 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.) $\bigcirc \ v_2=10$, $v_2^{corrected}=7.5$ $igcup v_2=7.5$, $v_2^{corrected}=10$ ✓ Correct Which of these is NOT a good learning rate decay scheme? Here, t is the epoch number. 1/1 point $igotimes lpha = e^t lpha_0$ $\bigcirc \; \; lpha = 0.95^t lpha_0$ $\bigcirc \ \ lpha = rac{1}{1+2*t}lpha_0$ $\bigcirc \ \ lpha = rac{1}{\sqrt{t}}lpha_0$ ✓ Correct 7. You use an exponentially weighted average on the London temperature dataset. You use the 1 / 1 point following to track the temperature: $v_t = eta v_{t-1} + (1-eta) heta_t$. The red line below was computed using perature days Increasing β will shift the red line slightly to the right. ✓ Correct True, remember that the red line corresponds to eta=0.9. In lecture we had a green line \$\$\beta = 0.98) that is slightly shifted to the right. Decreasing eta will create more oscillation within the red line. ✓ Correct True, remember that the red line corresponds to eta=0.9. In lecture we had a yellow line \$\$\beta = 0.98 that had a lot of oscillations. Increasing β will create more oscillations within the red line. Consider this figure: 1/1 point 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 with momentum (small β). (2) is gradient descent. (3) is gradient descent with momentum (large β) (1) is gradient descent with momentum (small β), (2) is gradient descent with momentum (small β), (3) is gradient descent (1) is gradient descent. (2) is gradient descent with momentum (small β). (3) is gradient descent with momentum (large β) (1) is gradient descent. (2) is gradient descent with momentum (large eta) . (3) is gradient descent with momentum (small β) Suppose batch gradient descent in a deep network is taking excessively long to find a value of the 1 / 1 point 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 ${\cal J}$? (Check all that apply) Try mini-batch gradient descent Correct Try initializing all the weights to zero Try better random initialization for the weights Try tuning the learning rate lphaCorrect Try using Adam ✓ Correct 10. 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 ($eta_1=0.9,eta_2=0.999$ $\epsilon = 10^{-8}$

Adam combines the advantages of RMSProp and momentum

Adam should be used with batch gradient computations, not with mini-batches.

The learning rate hyperparameter lpha in Adam usually needs to be tuned.

Due Jan 13, 1:59 AM CST