On-policy Prediction with Approximation Practice Quiz • 30 min GRADE Congratulations! You passed! 100% **Keep Learning** TO PASS 80% or higher **On-policy Prediction with Approximation TOTAL POINTS 12** 1. Which of the following statements is true about function approximation in reinforcement learning? 1 / 1 point (Select all that apply) It allows faster training by generalizing between states. Correct Correct! Function approximation allows the agent to generalize to unseen but similar states, and can learn the value function more quickly. Furthermore, in continuous state/action spaces the agent may never see the same state twice and we need such generalization to accurately estimate the values. It can help the agent achieve good generalization with good discrimination, so that it learns faster and represent the values quite accurately. Correct Correct. Recall the 2D plot of generalization and discrimination. Tabular methods discriminate between different states perfectly but with no generalization. Alternatively, one could treat all states as the same, with each update generalizing to all states but with no discrimination. Ideal function approximation methods achieves both good generalization and good discrimination. We only use function approximation because we have to for large or continuous state spaces. We would use tabular methods if we could, and learn an independent value per state. It can be more memory efficient. Correct Correct! We cannot enumerate and store all states in a table for large or continuous state spaces. By using function approximation, we can use fewer parameters to represent the value function. 2. We learned how value function estimation can be framed as supervised learning. But not all 1 / 1 point supervised learning methods are suitable. What are some key differences in reinforcement learning that can make it hard to apply standard supervised learning methods? Data is available as a fixed batch. Data is temporally correlated in reinforcement learning. Correct Correct! When using bootstrapping methods like TD, the target labels change. Correct Correct. Targets depend on our own estimates, and these estimates change as learning progresses. Reinforcement learning is usually done in an online setting where the full dataset is not fixed and unavailable from the beginning. Correct Correct! 3. Value Prediction (or Policy Evaluation) with Function Approximation can be viewed as supervised 1 / 1 point learning mainly because \_\_\_\_\_. [choose the most appropriate completion of the proceeding statement] Each state and its target estimate (used in the Monte Carlo update, TD(0) update, and DP update) can be seen as input-output training examples to estimate a continuous function. We can learn the value function by training with batches of data obtained from the agent's interaction with the world. We use stochastic gradient descent to learn the value function. Correct They can be seen as an (input, output) training example with  $(S_t,G_t)$  for Monte Carlo update,  $(S_t, R_{t+1} + \gamma V_\pi(S_{t+1}))$  for TD(0) update, and  $(s, E_\pi(R_{t+1} + \gamma V_\pi(S_{t+1})|S_t = s))$ for DP update. Each of these updates makes the output of value function of  $S_t$  more like the target value. 4. Which of the following is true about using Mean Squared Value Error ( 1 / 1 point  $V\bar{E} = \sum \mu(s)[v_{\pi}(s) - \hat{v}(s, w)]^2$ ) as the prediction objective? (Select all that apply) The agent can get zero MSVE when using a linear representation that cannot represent the true values. Gradient Monte Carlo with linear function approximation converges to the global optimum of this objective, if the step size is reduced over time. Correct Correct. There are stronger theoretical guarantees with linear function approximation than with non-linear function approximation. This objective makes it explicit how we should trade-off accuracy of the value estimates across states, using the weighting  $\mu$ . Correct Correct.  $\mu(s)$  is a weighting of how much we care about the error in state s, and we usually choose  $\mu(s)$  to be the fraction of time we spend in state s. The agent can get zero MSVE when using a tabular representation that can represent the true values. Correct Correct. In fact, in the tabular setting, we did not define an objective because we did not need to. With a table of values, we can represent the true value function exactly. So we do not need an objective to help specify how to trade-off accuracy. 5. Which of the following is true about  $\mu(S)$  in Mean Squared Value Error? (Select all that apply) 1 / 1 point It is a probability distribution. Correct Correct. If the policy is uniformly random,  $\mu(S)$  would have the same value for all states. It has higher values for states that are visited more often. Correct Correct. It serves as a weighting to minimize the error more in states that we care about. Correct Correct. 6. If we are given the true value function  $v_\pi(S_t)$ , the stochastic gradient descent update would be as 1 / 1 point follows. Fill in the blanks (A), (B), (C) and (D) with correct terms. (Select all correct answers)  $\mathbf{w_{t+1}} \doteq \mathbf{w_t} (A) \frac{1}{2} \alpha \nabla [(C) - (D)]^2$  $= \mathbf{w_t} (B) \alpha [(C) - (D)] \nabla \hat{\mathbf{v}}(S_t, \mathbf{w_t})$  $(\alpha > 0)$ -, +,  $v_{\pi}(S_t)$ ,  $\mathring{v}(S_t, \mathbf{w_t})$ Correct Correct! stochastic gradient descent makes update to  $\mathbf{w}_t$  proportional to the <u>negative</u> gradient of the squared error.  $-, -, \hat{\mathcal{V}}(S_t, \mathbf{w_t}), \nu_{\pi}(S_t)$ Correct Correct! stochastic gradient descent makes update to  $\mathbf{w}_t$  proportional to the <u>negative</u> gradient of the squared error.  $\square$  +, -,  $\nu_{\pi}(S_t)$ ,  $\mathring{\nu}(S_t, \mathbf{w_t})$  $\square$  +, +,  $\mathring{v}(S_t, \mathbf{w_t}), v_{\pi}(S_t)$ 7. In a Monte Carlo Update with function approximation, we do stochastic gradient descent using the 1 / 1 point following gradient:  $\nabla [G_t - \hat{\mathbf{v}}(s, \mathbf{w})]^2 = 2[G_t - \hat{\mathbf{v}}(s, \mathbf{w})] \nabla (-\hat{\mathbf{v}}(S_t, \mathbf{w}_t))$  $= (-1) * 2[G_t - \hat{v}(s, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w}_t)$ But the actual Monte Carlo Update rule is the following:  $\mathbf{w_{t+1}} = \mathbf{w_t} + \alpha [G_t - \hat{v}(S_t, \mathbf{w}_t)] \nabla \hat{v}(S_t, \mathbf{w}_t), \qquad (\alpha > 0)$ Where did the constant -1 and 2 go when  $\alpha$  is positive? (Choose all that apply) We assume that the 2 is included in  $\nabla \hat{v}(S_t, \mathbf{w}_t)$ . We are performing gradient ascent, so we subtract the gradient from the weights, negating -1. We are performing gradient descent, so we subtract the gradient from the weights, negating -1. Correct Correct. We assume that the 2 is included in the step-size. Correct Correct. It is equivalent to use  $\alpha$  or  $2\alpha$ , because we select  $\alpha$ . If we want to use an  $\alpha$  of 0.1 for the gradient with a 2 in front, then it is equivalent to use an lpha of 0.2 without a 2 in front of the gradient. 8. When using stochastic gradient descent for learning the value function, why do we only make a 1 / 1 point small update towards minimizing the error instead of fully minimizing the error at each encountered state? Because the target value may not be accurate initially for both TD(0) and Monte Carlo method. Because small updates guarantee we can slowly reduce approximation error to zero for all states. Because we want to minimize approximation error for all states, proportionally to  $\mu$ . Correct Correct! With function approximation, the agents have limited capacity and minimizing the approximation error for one state invariably increases the error for other states. We want to make small updates so that the error is reduced across states, proportionally to the weighting  $\mu$ . 9. The general stochastic gradient descent update rule for state-value prediction is as follows: 1 / 1 point  $\mathbf{w_{t+1}} \doteq \mathbf{w_t} + \alpha [U_t - \hat{v}(S_t, \mathbf{w_t})] \nabla \hat{v}(S_t, \mathbf{w_t})$ For what values of  $U_t$  would this be a semi-gradient method?  $\bigcirc$   $R_{t+1} + \hat{v}(S_{t+1}, w_t)$  $\bigcap R_{t+1} + R_{t+2} + ... + R_T$  $\bigcirc v_{\pi}(S_t)$  $\bigcirc G_t$ Correct Correct. This is the typical TD(0) bootstrapping target, which depends on the current weight vector  $\mathbf{w_t}$ . It will not produce a true gradient estimate, because its expected value is not equal to true  $v_{\pi}$ . 10. Which of the following statements is true about state-value prediction using stochastic gradient 1 / 1 point descent?  $\mathbf{w_{t+1}} \doteq \mathbf{w_t} + \alpha [U_t - \hat{v}(S_t, \mathbf{w_t})] \nabla \hat{v}(S_t, \mathbf{w_t})$ (Select all that apply) Stochastic gradient descent updates with Monte Carlo targets always reduce the Mean Squared Value Error at each step. Semi-gradient TD(0) methods typically learns faster than gradient Monte Carlo methods. Correct Correct! Similar to the tabular case, Semi-gradient TD(0) methods learn faster than gradient Monte Carlo methods. Using the Monte Carlo return as target, and under appropriate stochastic approximation conditions, the value function will converge to a local optimum of the Mean Squared Value Error. Correct Correct! Monte Carlo return ( $G_t$ ) is an unbiased estimate of  $v_{\pi}(S_t)$ . It converges to a stationary point, which under mild conditions, will be a local optimum of the MSVE. Using the Monte Carlo return or true value function as target results in an unbiased update. Correct True. The stochastic update with either target is an unbiased estimate of the gradient of the MSVE. When using  $U_t = R_{t+1} + \hat{v}(S_{t+1}, \mathbf{w_t})$ , the weight update is not using the true gradient of the TD error. Correct Correct! When computing the gradient of the TD error, we do not consider the effect of changing the weight vector  $\mathbf{w_t}$  in the bootstrapped target  $U_t$ . 11. Which of the following is true about the TD fixed point? 1 / 1 point (Select all correct answers) The weight vector corresponding to the TD fixed point is a local minimum of the Mean Squared Value Error. Semi-gradient TD(0) with linear function approximation converges to the TD fixed point. Correct Correct! This is the definition of TD fixed point. At the TD fixed point, the mean squared value error is not larger than  $\frac{1}{1-\gamma}$  times the mean squared value error of the global optimum, assuming the same linear function approximation. Correct Correct! See Equation (9.14) from the textbook. The weight vector corresponding to the TD fixed point is the global minimum of the Mean Squared Value Error. 12. Which of the following is true about Linear Function Approximation, for estimating state-values? 1 / 1 point (Select all that apply) State aggregation is one way to generate features for linear function approximation. Correct Correct. Features are often called basis functions because every approximate value function we consider can be written as a linear combination of these features. Correct Correct. The size of the feature vector is not necessarily equal to the size of the weight vector. The gradient of the approximate value function  $\hat{v}(s, \mathbf{w})$  with respect to  $\mathbf{w}$  is just the feature vector. Correct Correct. In linear function approximation, the value function is a linear combination of the weight vector and the feature vector.  $\hat{v}(s, \mathbf{w}) = \mathbf{w}^T \mathbf{x}(s)$ . By taking the gradient with respect to w, the gradient is the feature vector x(s).