



✓ Поздравляем! Вы прошли тест!

для успешного прохождения 80% или выше

Продолжить обучение

ОЦЕНКА
94.44%

Policy Gradient Methods

ОБЩИЙ БАЛЛ 18

1. Which of the following is true about policy-based methods? (Select all that apply)

1 / 1 балл

Policy-based methods can be applied to continuous action space domains.

✓ Correct

Correct. By parameterizing a policy to represent a probability distribution such as Gaussian, it can be applied to continuous action space domains.

Policy-based methods allow smooth improvement in the policy without drastic changes.

✓ Correct

Correct. As the policy parameters change the action probabilities change smoothly, but with value-based methods a small change in action-value function can drastically change the action probabilities.

Policy-based methods are useful in problems where the policy is easier to approximate than action-value functions.

✓ Correct

Correct. For example in the Mountain Car problem a good policy is easy to represent whereas the value function is complex.

Policy-based methods can learn an optimal policy that is stochastic.

✓ Correct

Correct. It can learn a stochastic optimal policy, such as the soft-max in action preferences.

2. Which of the following statements about parameterized policies are true? (Select all that apply)

1 / 1 балл

The policy must be approximated using linear function approximation.

For each state, the sum of all the action probabilities must equal to one.

✓ Correct

Correct! This condition is necessary for the function to be a valid probability distribution.

The function used for representing the policy must be a softmax function.

The probability of selecting any action must be greater than or equal to zero.

✓ Correct

Correct! This is one of the conditions for a valid probability distribution.

3. Assume you're given the following preferences $h_1 = 44$, $h_2 = 42$, and $h_3 = 38$, corresponding to three different actions (a_1, a_2, a_3), respectively. Under a softmax policy, what is the probability of choosing a_2 , rounded to three decimal numbers?

1 / 1 балл

0.119

0.42

0.879

0.002

✓ Correct

Correct!

4. Which of the following is true about softmax policy? (Select all that apply)

1 / 1 балл

It cannot represent an optimal policy that is stochastic, because it reaches a deterministic policy as one action preference dominates others.

It can be parameterized by any function approximator as long as it can output scalar values for each available action, to form a softmax policy.

Correct
Correct. It can use any function approximation from deep artificial neural networks to simple linear features.

It is used to represent a policy in discrete action spaces.

Correct
Correct!

Similar to epsilon-greedy policy, softmax policy cannot approach a deterministic policy.

5. What are the differences between using softmax policy over action-values and using softmax policy over action-preferences? (Select all that apply)

1 / 1 балл

- When using softmax policy over action-values, even if the optimal policy is deterministic, the policy may never approach a deterministic policy.

Correct
Correct. The policy will always select proportional to exponentiated action-values.

When using softmax policy over action-values, assuming a tabular representation, the policy will converge to the optimal policy regardless of whether the optimal policy is stochastic or deterministic.

When using softmax policy over action-preferences, assuming a tabular representation, the policy will converge to the optimal policy regardless of whether the optimal policy is stochastic or deterministic.

Correct
Correct. Action-preferences does not approach specific values like action-values do. They can be driven to produce a stochastic policy or deterministic policy.

6. What is the following objective, and in which task formulation?

1 / 1 балл

$$r(\pi) = \sum_s \mu(s) \sum_a \pi(a|s, \theta) \sum_{s', r} p(s', r|s, a) r$$

- Average reward objective, continuing task
- Discounted return objective, continuing task
- Undiscounted return objective, episodic task
- Average reward objective, episodic task

Correct
Correct.

7. Which of the following is true about policy gradient methods? (Select all that apply)

1 / 1 балл

- If we have access to the true value function v_π , we can perform unbiased stochastic gradient updates using the result from the Policy Gradient Theorem.

Correct
Correct. We derived this stochastic update by multiplying and dividing by $\pi(A|S)$.

- The policy gradient theorem provides a form for the policy gradient that does not contain the gradient of the state distribution $\lambda\mu$, which is hard to estimate.

Correct
Correct.

- Policy gradient methods use generalized policy iteration to learn policies directly.

- Policy gradient methods do gradient ascent on the policy objective.

Correct
Correct. Policy gradient methods maximize the policy objective, and hence perform gradient ascent.

8. The following equation is the outcome of the policy gradient theorem. Which of the following is true about the policy gradient theorem? (Select all that apply)

1 / 1 балл

$$\nabla r(\pi) = \sum_s \mu(s) \sum_a \nabla \pi(a|s, \theta) q_\pi(s, a)$$

- This expression can be converted into the following expectation over π :

$$\mathbb{E}_\pi[\nabla \ln \pi(A|S, \theta) q_\pi(S, A)]$$

Correct
Correct. In fact, this expression is normally used to perform stochastic gradient updates.

- We do not need to compute the gradient of the state distribution μ .

✓ Correct
Correct.

- This expression can be converted into:

$$\mathbb{E}_\pi[\sum_a \nabla \pi(a|S, \theta) q_\pi(S, a)]$$

In discrete action space, by approximating q_{pi} we could also use this gradient to update the policy.

✓ Correct

Correct. The expression contains sum over actions, which can be computed for discrete actions. In the textbook, this is also called the all-actions method.

- The true action value q_π can be approximated in many ways, for example using TD algorithms.

✓ Correct
Correct.

9. Which of the following statements is true? (Select all that apply)

1 / 16 ann

- TD methods do not have a role when estimating the policy directly.
 To update the actor in Actor-Critic, we can use TD error in place of q_π in the Policy Gradient Theorem.

✓ Correct

Correct. This is equivalent to using one-step state value and subtracting a current state value baseline.

- The Actor-Critic algorithm consists of two parts: a parameterized policy — the actor — and a value function — the critic.

✓ Correct
Correct.

- Subtracting a baseline in the policy gradient update tends to reduce the variance of the update, which results in faster learning.

✓ Correct
Correct.

10. To train the critic, we must use the average reward version of semi-gradient TD(0).

1 / 16 ann

- False
 True

✓ Correct

Correct. We can use any state-value learning algorithm.

11. Question 11 ~ 13: Consider the following state features and parameters θ for three different actions (red, green, and blue):

1 / 16 ann

$$\mathbf{X}(s) = \begin{bmatrix} 0.1 \\ 0.3 \\ 0.6 \end{bmatrix} \quad \theta = \begin{bmatrix} 45 \\ 73 \\ 21 \\ 120 \\ 120 \\ -10 \\ -100 \\ 200 \\ -25 \end{bmatrix} \left\{ \begin{array}{l} a_0 \\ a_1 \\ a_2 \end{array} \right\}$$

Compute the action preferences for each of the three different actions using linear function

approximation and stacked features for the action preferences.

What is the action preference of a_0 (red)?

- 39
- 37
- 35
- 33

✓ Correct
Correct.

12. What is the action preference of a_1 (green)?

- 40
- 35
- 32
- 42

✓ Correct
Correct.

13. What is the action preference of a_2 (blue)?

- 37
- 42
- 35
- 39

✓ Correct
Correct.

14. Which of the following statements are true about the Actor-Critic algorithm with softmax policies?
(Choose all that apply)

- The actor and the critic share the same set of parameters.
- The preferences must be approximated using linear function approximation.
- The learning rate parameter of the actor and the critic can be different.

✓ Correct
Correct! In practice, it is preferable to have a slower learning rate for the actor so that the critic can accurately critique the policy.

- Since the policy is written as a function of the current state, it is like having a different softmax distribution for each state.

✓ Correct
Correct!

15. We usually want the critic to update at a faster rate than the actor.

- True
- False

✓ Correct
Correct.

16. Which one is a reasonable parameterization for a Gaussian policy?

- μ : a linear function of parameters, σ : the exponential of a linear function of parameters.
- μ : the exponential of a linear function of parameters, σ : a linear function of parameters.
- μ : a linear function of parameters, σ : a linear function of parameters

! Incorrect
Incorrect. Remember that the parameter sigma must be positive. A linear function does not guarantee that this constraint will be met.

17. A Gaussian policy becomes deterministic in the limit $\sigma \rightarrow 0$.

- True
- False

1 / 1 bann

1 / 1 bann

1 / 1 bann

1 / 1 bann

0 / 1 bann

1 / 1 bann

 Correct

Correct: As σ approaches 0, the values of the Gaussian policy approach the mean of the policy in a given state.

1 / 16 ann

18. Which of the following is an advantage of Gaussian policy parameterization over discretizing the action space? (Select all that apply)

There might not be a straightforward way to choose a discrete set of actions.

 Correct

Correct! Selecting a discrete set of actions that results in good performance is problem dependent. Maybe we need hundreds of actions. Maybe it is state dependent!

Continuous actions also allow learning to generalize over actions.

 Correct

Correct!

Even if the true action set is discrete, but very large, it might be better to treat them as a continuous range.

 Correct

Correct!

Gaussian policies are differentiable, whereas policies over discretized actions are not.