# Actor Critic II

## CISC 7404 - Decision Making

Steven Morad

University of Macau

# Quiz

# Quiz

- After all students put away their computer/notes/phone I will hand out exams

# Quiz

- After all students put away their computer/notes/phone I will hand out exams
- If you have computer/notes/phone out after this point, it counts as cheating

# Quiz

- After all students put away their computer/notes/phone I will hand out exams
- If you have computer/notes/phone out after this point, it counts as cheating
- I will hand out exams face down, do not turn them over until I say so

# Quiz

- After all students put away their computer/notes/phone I will hand out exams
- If you have computer/notes/phone out after this point, it counts as cheating
- I will hand out exams face down, do not turn them over until I say so
- After turning them over, I will briefly explain each question

# Quiz

- After all students put away their computer/notes/phone I will hand out exams
- If you have computer/notes/phone out after this point, it counts as cheating
- I will hand out exams face down, do not turn them over until I say so
- After turning them over, I will briefly explain each question
- After my explanation, you will have 75 minutes to complete the exam

# Quiz

- After all students put away their computer/notes/phone I will hand out exams
- If you have computer/notes/phone out after this point, it counts as cheating
- I will hand out exams face down, do not turn them over until I say so
- After turning them over, I will briefly explain each question
- After my explanation, you will have 75 minutes to complete the exam
- After you are done, give me your exam and go relax outside, we resume class at 8:30

# Quiz

- There may or may not be different versions of the exam

# Quiz

- There may or may not be different versions of the exam
- If your exam has the answer for another version, it is cheating

# Quiz

- There may or may not be different versions of the exam
- If your exam has the answer for another version, it is cheating
- Instructions are in both english and chinese, english instructions take precedence

# Quiz

- There may or may not be different versions of the exam
- If your exam has the answer for another version, it is cheating
- Instructions are in both english and chinese, english instructions take precedence
- Good luck!

# Quiz

- 在所有学生收起电脑/笔记/手机后,我会分发试卷。
- 如果在此之后仍有电脑/笔记/手机未收,将视为作弊。
- 试卷会背面朝下发下,在我宣布开始前请勿翻面。
- 试卷翻面后,我会简要说明每道题的注意事项。
- 说明结束后,你们有 75 分钟完成考试。
- 交卷后请到教室外休息,8:30 恢复上课。
- 试卷可能存在不同版本,细节略有差异。
- 若你的试卷上出现其他版本的答案,将被判定为作弊。
- 试卷说明为中英双语,若内容冲突以英文为准。
- 祝各位考试顺利!

# Admin

# Admin

After today, you will know all the foundations of decision making

# Admin

After today, you will know all the foundations of decision making

You should be able to understand research papers on:

- MDPs, bandits, RL algorithms, even some control problems!

# Admin

After today, you will know all the foundations of decision making

You should be able to understand research papers on:
- MDPs, bandits, RL algorithms, even some control problems!

Next week, we begin to investigate interesting topics of decision making

# Admin

After today, you will know all the foundations of decision making

You should be able to understand research papers on:
- MDPs, bandits, RL algorithms, even some control problems!

Next week, we begin to investigate interesting topics of decision making

Right now, my plan is:
- Imitation Learning
- Offline RL
- Memory and POMDPs
- Large Language Models

# Admin

After today, you will know all the foundations of decision making

You should be able to understand research papers on:
- MDPs, bandits, RL algorithms, even some control problems!

Next week, we begin to investigate interesting topics of decision making

Right now, my plan is:
- Imitation Learning
- Offline RL
- Memory and POMDPs
- Large Language Models

**Question:** Should we replace a topic with something else?

# Admin

- Imitation learning
  - ‣ Sometimes, designing a reward function is hard
  - ‣ It is easier to demonstrate desired behavior to agents
    - – Instead of reward for surgery, do what human does
    - – Instead of reward for self driving, do what human does
  - ‣ With imitation learning, can learn behaviors without rewards
  - ‣ Closer to supervised learning, easier to train
  - ‣ Policies are not better than dataset/humans

# Admin

- Offline RL
  - ‣ RL without exploration
  - ‣ How can we learn policies from a fixed dataset?
    - – Learn surgery from surgical videos (no need to kill patients)
    - – Learn driving from Xiaomi driving dataset (no need to crash cars)
  - ‣ Unlike imitation learning, can do **better** than dataset
  - ‣ Very new topic (2-3 years old)
    - – Does not work very well (yet)

# Admin

- Memory and POMDPs (my research focus)
  - ▸ So far, we focused on video games
    - – MDP
  - ▸ Many interesting problems are not Markov
    - – Think of robot with camera, not Markov
    - – Almost every task has sensor noise, not Markov
  - ▸ Can we extend RL to work for virtually any problem?
    - – Yes, requires long-term memory
    - – LSTM, transformer, etc
  - ▸ May also have time to introduce world models
    - – Dreamer, TD-MPC, etc

# Admin

- Large Language Models
  - ‣ Can train LLMs using unsupervised learning
    - – They only learn to predict next word
  - ‣ We use RL to teach them to interact with humans
    - – Apply policy gradient to textual MDP
    - – DeepSeek math/GRPO
    - – RL-adjacent methods (DPO)

# Admin

Also a possibility to split lecture:

- E.g., 1 hour imitation learning, 1 hour something new

# Admin

Also a possibility to split lecture:

- E.g., 1 hour imitation learning, 1 hour something new

**Question:** Any topic sound boring?

# Admin

Also a possibility to split lecture:

- E.g., 1 hour imitation learning, 1 hour something new

**Question:** Any topic sound boring?

**Question:** Any suggestions for other topics?

# Admin

Also a possibility to split lecture:

- E.g., 1 hour imitation learning, 1 hour something new

**Question:** Any topic sound boring?

**Question:** Any suggestions for other topics?

- Maybe just focus on a specific paper?

# Admin

Also a possibility to split lecture:

- E.g., 1 hour imitation learning, 1 hour something new

**Question:** Any topic sound boring?

**Question:** Any suggestions for other topics?

- Maybe just focus on a specific paper?

Alternative topics:

- Multi-agent RL
- Model-based RL and world-models
- Evolutionary algorithms

# Admin

Homework 2 progress

# Admin

Homework 2 progress

If you did not already start, you might be in trouble

# Admin

Homework 2 progress

If you did not already start, you might be in trouble

Experiments take a long time, start as soon as possible

# Admin

Homework 2 progress

If you did not already start, you might be in trouble

Experiments take a long time, start as soon as possible

Harder and requires more debugging than `FrozenLake` assignment

# Admin

Those using Tencent AI Arena (Honor of Kings):

# Admin

Those using Tencent AI Arena (Honor of Kings):

- Backup project should not rely on Honor of Kings

# Admin

Those using Tencent AI Arena (Honor of Kings):

- Backup project should not rely on Honor of Kings
- **Make sure to save your code regularly**

# Admin

Those using Tencent AI Arena (Honor of Kings):

- Backup project should not rely on Honor of Kings
- **Make sure to save your code regularly**
  - ‣ Everything stored on Tencent VMs, not sure how safe code is

# Admin

Those using Tencent AI Arena (Honor of Kings):

- Backup project should not rely on Honor of Kings
- **Make sure to save your code regularly**
  - ‣ Everything stored on Tencent VMs, not sure how safe code is
- There is already a PPO baseline, cannot use PPO

# Admin

Those using Tencent AI Arena (Honor of Kings):

- Backup project should not rely on Honor of Kings
- **Make sure to save your code regularly**
  - ‣ Everything stored on Tencent VMs, not sure how safe code is
- There is already a PPO baseline, cannot use PPO
  - ‣ Instead, consider DDPG, SAC, TRPO, DQN variant, etc

# Admin

Those using Tencent AI Arena (Honor of Kings):

- Backup project should not rely on Honor of Kings
- **Make sure to save your code regularly**
  - ‣ Everything stored on Tencent VMs, not sure how safe code is
- There is already a PPO baseline, cannot use PPO
  - ‣ Instead, consider DDPG, SAC, TRPO, DQN variant, etc
- Cannot install new python libraries (Tencent security issue)

# Admin

Those using Tencent AI Arena (Honor of Kings):

- Backup project should not rely on Honor of Kings
- **Make sure to save your code regularly**
  - ‣ Everything stored on Tencent VMs, not sure how safe code is
- There is already a PPO baseline, cannot use PPO
  - ‣ Instead, consider DDPG, SAC, TRPO, DQN variant, etc
- Cannot install new python libraries (Tencent security issue)
  - ‣ No `jax`, must use `torch`
  - ‣ You must learn Tencent's strange callback system
    - – Prevents copy/pasting, so `torch` is ok
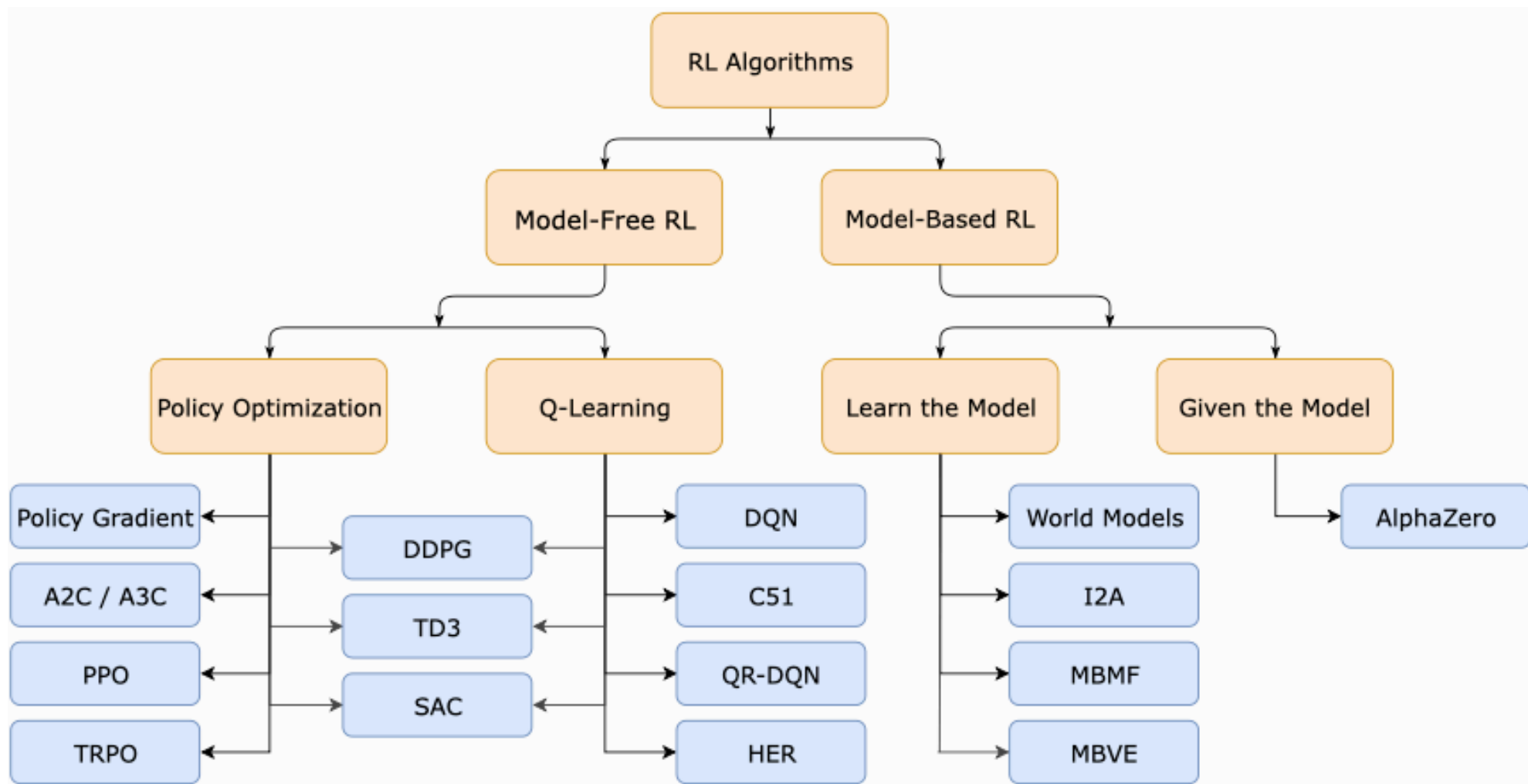
# Review

# Actor Critic

# Actor Critic

Alternative descriptions of actor critic algorithms

https://lilianweng.github.io/posts/2018-04-08-policy-gradient/

# Actor Critic

# Actor Critic

There are two approaches to actor critic

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

PG with $V$ instead of MC

# Actor Critic

There are two approaches to actor critic

1.  **Policy gradient based:**

PG with $V$ instead of MC

- A2C

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

PG with $V$ instead of MC

- A2C
- TRPO

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

PG with $V$ instead of MC

- A2C
- TRPO
- PPO

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

PG with $V$ instead of MC

- A2C
- TRPO
- PPO
- IMPALA

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

PG with $V$ instead of MC

- A2C
- TRPO
- PPO
- IMPALA
- ACKTR

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

PG with $V$ instead of MC

- A2C
- TRPO
- PPO
- IMPALA
- ACKTR
- ACER

# Actor Critic

There are two approaches to actor critic

1.  **Policy gradient based:**

PG with $V$ instead of MC

- A2C
- TRPO
- PPO
- IMPALA
- ACKTR
- ACER
- PPG

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**        2. **Q learning based:**

PG with $V$ instead of MC
- A2C
- TRPO
- PPO
- IMPALA
- ACKTR
- ACER
- PPG

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

   PG with $V$ instead of MC
   - A2C
   - TRPO
   - PPO
   - IMPALA
   - ACKTR
   - ACER
   - PPG

2. **Q learning based:**

   Learn $Q$ for a specific policy

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

   PG with $V$ instead of MC
   - A2C
   - TRPO
   - PPO
   - IMPALA
   - ACKTR
   - ACER
   - PPG

2. **Q learning based:**

   Learn $Q$ for a specific policy
   - DDPG

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

PG with $V$ instead of MC
- A2C
- TRPO
- PPO
- IMPALA
- ACKTR
- ACER
- PPG

2. **Q learning based:**

Learn $Q$ for a specific policy
- DDPG
- TD3

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

PG with $V$ instead of MC

- A2C
- TRPO
- PPO
- IMPALA
- ACKTR
- ACER
- PPG

2. **Q learning based:**

Learn $Q$ for a specific policy

- DDPG
- TD3
- SAC

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

   PG with $V$ instead of MC
   - A2C
   - TRPO
   - PPO
   - IMPALA
   - ACKTR
   - ACER
   - PPG

2. **Q learning based:**

   Learn $Q$ for a specific policy
   - DDPG
   - TD3
   - SAC
   - Q-Prop

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

   PG with $V$ instead of MC
   - A2C
   - TRPO
   - PPO
   - IMPALA
   - ACKTR
   - ACER
   - PPG

2. **Q learning based:**

   Learn $Q$ for a specific policy
   - DDPG
   - TD3
   - SAC
   - Q-Prop
   - QT-Opt

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

   PG with $V$ instead of MC
   - A2C
   - TRPO
   - PPO
   - IMPALA
   - ACKTR
   - ACER
   - PPG

2. **Q learning based:**

   Learn $Q$ for a specific policy
   - DDPG
   - TD3
   - SAC
   - Q-Prop
   - QT-Opt
   - NAF

# Actor Critic

There are two approaches to actor critic

1. **Policy gradient based:**

   PG with $V$ instead of MC
   - A2C
   - TRPO
   - PPO
   - IMPALA
   - ACKTR
   - ACER
   - PPG

2. **Q learning based:**

   Learn $Q$ for a specific policy
   - DDPG
   - TD3
   - SAC
   - Q-Prop
   - QT-Opt
   - NAF
   - MPO

# Deterministic Policy Gradient

# Deterministic Policy Gradient

Question: Why did we introduce policy gradient methods?

# Deterministic Policy Gradient

**Question:** Why did we introduce policy gradient methods?

# Deterministic Policy Gradient

# Deterministic Policy Gradient



**Question:** Why did Q learning fail BenBen?

# Deterministic Policy Gradient



**Question:** Why did Q learning fail BenBen?

$$A = [0, 2\pi]^{12}$$

# Deterministic Policy Gradient

**Question:** Why did Q learning fail BenBen?

$$A = [0, 2\pi]^{12}$$

$$\pi(a_t \mid s_t; \theta_\pi) = \begin{cases} 1 \text{ if } a_t = \arg\max_{a_t \in A} Q(s_t, a_t, \theta_\pi) \\ 0 \text{ otherwise} \end{cases}$$

# Deterministic Policy Gradient

**Question:** Why did Q learning fail BenBen?

$$A = [0, 2\pi]^{12}$$

$$\pi(a_t \mid s_t; \theta_\pi) = \begin{cases} 1 \text{ if } a_t = \arg \max_{a_t \in A} Q(s_t, a_t, \theta_\pi) \\ 0 \text{ otherwise} \end{cases}$$

Infinitely many $a_t$ – compute $Q$ for each and take arg max over all

# Deterministic Policy Gradient

$$\pi(a_t \mid s_t; \theta_\pi) = \begin{cases} 1 \text{ if } a_t = \arg \max_{a_t \in A} Q(s_t, a_t, \theta_\pi) \\ 0 \text{ otherwise} \end{cases}$$

# Deterministic Policy Gradient

$$\pi(a_t \mid s_t; \theta_\pi) = \begin{cases} 1 \text{ if } a_t = \arg\max_{a_t \in A} Q(s_t, a_t, \theta_\pi) \\ 0 \text{ otherwise} \end{cases}$$

The greedy policy cannot work with continuous action spaces

# Deterministic Policy Gradient

$$\pi(a_t \mid s_t; \theta_\pi) = \begin{cases} 1 \text{ if } a_t = \arg\max_{a_t \in A} Q(s_t, a_t, \theta_\pi) \\ 0 \text{ otherwise} \end{cases}$$

The greedy policy cannot work with continuous action spaces

Can we replace arg max with a different deterministic function?

# Deterministic Policy Gradient

$$\pi(a_t \mid s_t; \theta_\pi) = \begin{cases} 1 \text{ if } a_t = \arg \max_{a_t \in A} Q(s_t, a_t, \theta_\pi) \\ 0 \text{ otherwise} \end{cases}$$

The greedy policy cannot work with continuous action spaces

Can we replace arg max with a different deterministic function?

Remember that the Q function and greedy policy are different

# Deterministic Policy Gradient

$$\pi(a_t \mid s_t; \theta_\pi) = \begin{cases} 1 \text{ if } a_t = \arg\max_{a_t \in A} Q(s_t, a_t, \theta_\pi) \\ 0 \text{ otherwise} \end{cases}$$

The greedy policy cannot work with continuous action spaces

Can we replace arg max with a different deterministic function?

Remember that the Q function and greedy policy are different

We can learn Q for **any** policy, does not need to be greedy policy

# Deterministic Policy Gradient

$$\pi(a_t \mid s_t; \theta_\pi) = \begin{cases} 1 \text{ if } a_t = \arg \max_{a_t \in A} Q(s_t, a_t, \theta_\pi) \\ 0 \text{ otherwise} \end{cases}$$

The greedy policy cannot work with continuous action spaces

Can we replace arg max with a different deterministic function?

Remember that the Q function and greedy policy are different

We can learn Q for **any** policy, does not need to be greedy policy

Let us quickly review the Q function and value function

# Deterministic Policy Gradient

Start with general form of Temporal Difference Q function

# Deterministic Policy Gradient

Start with general form of Temporal Difference Q function

$$Q(s_0, a_0, \theta_\pi) = \underbrace{\mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0, a_0\big]}_{\text{Reward for taking } a_0} + \gamma \underbrace{V(s_1, \theta_\pi)}_{\mathcal{G} \text{ following } \theta_\pi}$$

# Deterministic Policy Gradient

Start with general form of Temporal Difference Q function

$$Q(s_0, a_0, \theta_\pi) = \underbrace{\mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0, a_0\big]}_{\text{Reward for taking } a_0} + \gamma \underbrace{V(s_1, \theta_\pi)}_{\mathcal{G} \text{ following } \theta_\pi}$$

We can replace $V$ with $Q$ if $a \sim \pi(\cdot \mid s_1; \theta_\pi)$

# Deterministic Policy Gradient

Start with general form of Temporal Difference Q function

$$Q(s_0, a_0, \theta_\pi) = \underbrace{\mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0]}_{\text{Reward for taking } a_0} + \gamma \underbrace{V(s_1, \theta_\pi)}_{\mathcal{G} \text{ following } \theta_\pi}$$

We can replace $V$ with $Q$ if $a \sim \pi(\cdot \mid s_1; \theta_\pi)$

$$Q(s_0, a_0, \theta_\pi) = \underbrace{\mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0]}_{\text{Reward for taking } a_0} + \gamma \underbrace{Q(s_1, a, \theta_\pi)}_{\mathcal{G} \text{ following } \theta_\pi}$$

# Deterministic Policy Gradient

Start with general form of Temporal Difference Q function

$$Q(s_0, a_0, \theta_\pi) = \underbrace{\mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0]}_{\text{Reward for taking } a_0} + \gamma \underbrace{V(s_1, \theta_\pi)}_{\mathcal{G} \text{ following } \theta_\pi}$$

We can replace $V$ with $Q$ if $a \sim \pi(\cdot \mid s_1; \theta_\pi)$

$$Q(s_0, a_0, \theta_\pi) = \underbrace{\mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0]}_{\text{Reward for taking } a_0} + \gamma \underbrace{Q(s_1, a, \theta_\pi)}_{\mathcal{G} \text{ following } \theta_\pi}$$

For the greedy policy, we can reduce $Q$ further

# Deterministic Policy Gradient

Start with general form of Temporal Difference Q function

$$Q(s_0, a_0, \theta_\pi) = \underbrace{\mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0]}_{\text{Reward for taking } a_0} + \gamma \underbrace{V(s_1, \theta_\pi)}_{\mathcal{G} \text{ following } \theta_\pi}$$

We can replace $V$ with $Q$ if $a \sim \pi(\cdot \mid s_1; \theta_\pi)$

$$Q(s_0, a_0, \theta_\pi) = \underbrace{\mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0]}_{\text{Reward for taking } a_0} + \gamma \underbrace{Q(s_1, a, \theta_\pi)}_{\mathcal{G} \text{ following } \theta_\pi}$$

For the greedy policy, we can reduce $Q$ further

$$Q(s_0, a_0, \theta_\pi) = \underbrace{\mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0]}_{\text{Reward for taking } a_0} + \underbrace{\max_{a \in A} \gamma Q(s_1, a, \theta_\pi)}_{\mathcal{G} \text{ following } \theta_\pi}$$

# Deterministic Policy Gradient

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \max_{a \in A} \gamma Q(s_1, a, \theta_\pi)$$

# Deterministic Policy Gradient

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \max_{a \in A} \gamma Q(s_1, a, \theta_\pi)$$

Cannot use the max Q function with BenBen

# Deterministic Policy Gradient

$$\cancel{Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \max_{a \in A} \gamma Q(s_1, a, \theta_\pi)}$$

Cannot use the max Q function with BenBen

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \gamma V(s_1, \theta_\pi)$$

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

# Deterministic Policy Gradient

$$\cancel{Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \max_{a \in A} \gamma Q(s_1, a, \theta_\pi)}$$

Cannot use the max Q function with BenBen

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \gamma V(s_1, \theta_\pi)$$

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** Can we use continuous $a$ with $Q$?

# Deterministic Policy Gradient

$$\cancel{Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \max_{a \in A} \gamma Q(s_1, a, \theta_\pi)}$$

Cannot use the max Q function with BenBen

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \gamma V(s_1, \theta_\pi)$$

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** Can we use continuous $a$ with $Q$? **Answer:** Yes

# Deterministic Policy Gradient

$$\sout{Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \max_{a \in A} \gamma Q(s_1, a, \theta_\pi)}$$

Cannot use the max Q function with BenBen

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \gamma V(s_1, \theta_\pi)$$

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** Can we use continuous $a$ with $Q$? **Answer:** Yes

Q function $\Rightarrow$ no problem, policy $\Rightarrow$ problem $\quad \textcolor{red}{\arg\max_{a \in A} Q(s, a, \theta_\pi)}$

# Deterministic Policy Gradient

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \max_{a \in A} \gamma Q(s_1, a, \theta_\pi)$$

Cannot use the max Q function with BenBen

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \gamma V(s_1, \theta_\pi)$$

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** Can we use continuous $a$ with $Q$? **Answer:** Yes

Q function $\Rightarrow$ no problem, policy $\Rightarrow$ problem $\quad \arg\max_{a \in A} Q(s, a, \theta_\pi)$

Greedy deterministic policy worked well for discrete actions

# Deterministic Policy Gradient

$$\cancel{Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \max_{a \in A} \gamma Q(s_1, a, \theta_\pi)}$$

Cannot use the max Q function with BenBen

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \gamma V(s_1, \theta_\pi)$$

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0, a_0] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** Can we use continuous $a$ with $Q$? **Answer:** Yes

Q function $\Rightarrow$ no problem, policy $\Rightarrow$ problem $\quad$ $\textcolor{red}{\arg\max_{a \in A} Q(s, a, \theta_\pi)}$

Greedy deterministic policy worked well for discrete actions

Can we learn a different deterministic policy for continuous actions?

# Deterministic Policy Gradient

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0, a_0\big] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** What method can we use to learn $\theta_\pi$?

# Deterministic Policy Gradient

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0, a_0\big] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** What method can we use to learn $\theta_\pi$?

**Answer:** Policy gradient

# Deterministic Policy Gradient

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0, a_0\big] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** What method can we use to learn $\theta_\pi$?

**Answer:** Policy gradient

Perhaps our solution will combine $Q$ with policy gradient

# Deterministic Policy Gradient

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0, a_0\big] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** What method can we use to learn $\theta_\pi$?

**Answer:** Policy gradient

Perhaps our solution will combine $Q$ with policy gradient

The trick is to consider a **deterministic** policy

# Deterministic Policy Gradient

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0, a_0\big] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** What method can we use to learn $\theta_\pi$?

**Answer:** Policy gradient

Perhaps our solution will combine $Q$ with policy gradient

The trick is to consider a **deterministic** policy

$$\cancel{a \sim \pi(\cdot \mid s; \theta_\pi)}$$

# Deterministic Policy Gradient

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0, a_0\big] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** What method can we use to learn $\theta_\pi$?

**Answer:** Policy gradient

Perhaps our solution will combine $Q$ with policy gradient

The trick is to consider a **deterministic** policy

$$\cancel{a \sim \pi(\cdot \mid s; \theta_\pi)} \qquad\qquad a = \mu(s, \theta_\mu)$$

# Deterministic Policy Gradient

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0, a_0\big] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** What method can we use to learn $\theta_\pi$?

**Answer:** Policy gradient

Perhaps our solution will combine $Q$ with policy gradient

The trick is to consider a **deterministic** policy

$$\cancel{a \sim \pi(\cdot \mid s; \theta_\pi)} \qquad\qquad a = \mu(s, \theta_\mu) \qquad\qquad \mu : S \times \Theta \mapsto A$$

# Deterministic Policy Gradient

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0, a_0\big] + \gamma Q(s_1, a, \theta_\pi); \quad a \sim \pi(\cdot \mid s_1; \theta_\pi)$$

**Question:** What method can we use to learn $\theta_\pi$?

**Answer:** Policy gradient

Perhaps our solution will combine $Q$ with policy gradient

The trick is to consider a **deterministic** policy

$$\cancel{a \sim \pi(\cdot \mid s; \theta_\pi)} \qquad a = \mu(s, \theta_\mu) \qquad \mu : S \times \Theta \mapsto A$$

$$Q(s_0, a_0, \theta_\pi) = \mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0, a_0\big] + \gamma Q\big(s_1, \mu(s_1, \theta_\mu), \theta_\pi\big)$$

# Deterministic Policy Gradient

So let us learn parameters $\theta_\mu$ for this deterministic function

# Deterministic Policy Gradient

So let us learn parameters $\theta_\mu$ for this deterministic function

Recall policy gradient, we find $\theta_\pi$ using gradient ascent

# Deterministic Policy Gradient

So let us learn parameters $\theta_\mu$ for this deterministic function

Recall policy gradient, we find $\theta_\pi$ using gradient ascent

$$\theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_{\pi,i}} \mathbb{E}\big[\mathcal{G}(\tau) \mid s_0; \theta_{\pi,i}\big]$$

# Deterministic Policy Gradient

So let us learn parameters $\theta_\mu$ for this deterministic function

Recall policy gradient, we find $\theta_\pi$ using gradient ascent

$$\theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_{\pi,i}} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_{\pi,i}\big]$$

We need to know $\nabla_{\theta_\pi} \mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi]$

# Deterministic Policy Gradient

So let us learn parameters $\theta_\mu$ for this deterministic function

Recall policy gradient, we find $\theta_\pi$ using gradient ascent

$$\theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_{\pi,i}} \mathbb{E}\big[\mathcal{G}(\tau) \mid s_0; \theta_{\pi,i}\big]$$

We need to know $\nabla_{\theta_\pi} \mathbb{E}[\mathcal{G}(\tau) \mid s_0; \theta_\pi]$

We found this for the stochastic policy

# Deterministic Policy Gradient

So let us learn parameters $\theta_\mu$ for this deterministic function

Recall policy gradient, we find $\theta_\pi$ using gradient ascent

$$\theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_{\pi,i}} \mathbb{E}\Big[\mathcal{G}(\tau) \mid s_0; \theta_{\pi,i}\Big]$$

We need to know $\nabla_{\theta_\pi} \mathbb{E}[\mathcal{G}(\tau) \mid s_0; \theta_\pi]$

We found this for the stochastic policy

Can we also find this for a deterministic policy?

# Deterministic Policy Gradient

So let us learn parameters $\theta_\mu$ for this deterministic function

Recall policy gradient, we find $\theta_\pi$ using gradient ascent

$$\theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_{\pi,i}} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_{\pi,i}\big]$$

We need to know $\nabla_{\theta_\pi} \mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi]$

We found this for the stochastic policy

Can we also find this for a deterministic policy?

$$\nabla_{\theta_\pi} \mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] \Rightarrow \nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$$

# Deterministic Policy Gradient

So let us learn parameters $\theta_\mu$ for this deterministic function

Recall policy gradient, we find $\theta_\pi$ using gradient ascent

$$\theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_{\pi,i}} \mathbb{E}\big[\mathcal{G}(\tau) \mid s_0; \theta_{\pi,i}\big]$$

We need to know $\nabla_{\theta_\pi} \mathbb{E}[\mathcal{G}(\tau) \mid s_0; \theta_\pi]$

We found this for the stochastic policy

Can we also find this for a deterministic policy?

$$\nabla_{\theta_\pi} \mathbb{E}[\mathcal{G}(\tau) \mid s_0; \theta_\pi] \Rightarrow \nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\tau) \mid s_0; \theta_\mu\big]$$

# Deterministic Policy Gradient

The expected return with a **stochastic** policy

# Deterministic Policy Gradient

The expected return with a **stochastic** policy

$$\mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \Pr(s_{n+1} \mid s_0; \theta_\pi)$$

# Deterministic Policy Gradient

The expected return with a **stochastic** policy

$$\mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \Pr(s_{n+1} \mid s_0; \theta_\pi)$$

# Deterministic Policy Gradient

The expected return with a **stochastic** policy

$$\mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \mathrm{Pr}(s_{n+1} \mid s_0; \theta_\pi)$$

$$\mathrm{Pr}(s_{n+1} \mid s_0; \theta_\pi) = \sum_{s_1,\ldots,s_n \in S} \prod_{t=0}^{n} \left( \sum_{a_t \in A} \mathrm{Tr}(s_{t+1} \mid s_t, a_t) \cdot \pi(a_t \mid s_t; \theta_\pi) \right)$$

# Deterministic Policy Gradient

The expected return with a **stochastic** policy

$$\mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \Pr(s_{n+1} \mid s_0; \theta_\pi)$$

$$\Pr(s_{n+1} \mid s_0; \theta_\pi) = \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \left( \sum_{a_t \in A} \mathrm{Tr}(s_{t+1} \mid s_t, a_t) \cdot \pi(a_t \mid s_t; \theta_\pi) \right)$$

# Deterministic Policy Gradient

The expected return with a **stochastic** policy

$$\mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \Pr(s_{n+1} \mid s_0; \theta_\pi)$$

$$\Pr(s_{n+1} \mid s_0; \theta_\pi) = \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \left( \sum_{a_t \in A} \mathrm{Tr}(s_{t+1} \mid s_t, a_t) \cdot \pi(a_t \mid s_t; \theta_\pi) \right)$$

# Deterministic Policy Gradient

The expected return with a **stochastic** policy

$$\mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \Pr(s_{n+1} \mid s_0; \theta_\pi)$$

$$\Pr(s_{n+1} \mid s_0; \theta_\pi) = \sum_{s_1,\dots,s_n \in S} \prod_{t=0}^{n} \left( \sum_{a_t \in A} \mathrm{Tr}(s_{t+1} \mid s_t, a_t) \cdot \pi(a_t \mid s_t; \theta_\pi) \right)$$

The state distribution with a **deterministic** policy

# Deterministic Policy Gradient

The expected return with a **stochastic** policy

$$\mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \Pr(s_{n+1} \mid s_0; \theta_\pi)$$

$$\Pr(s_{n+1} \mid s_0; \theta_\pi) = \sum_{s_1,\ldots,s_n \in S} \prod_{t=0}^{n} \left( \sum_{a_t \in A} \mathrm{Tr}(s_{t+1} \mid s_t, a_t) \cdot \pi(a_t \mid s_t; \theta_\pi) \right)$$

The state distribution with a **deterministic** policy

$$\Pr(s_{n+1} \mid s_0; \theta_\mu) = \sum_{s_1,\ldots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu))$$

# Deterministic Policy Gradient

Let us continue to derive policy gradient with a **deterministic** policy $\mu$

# Deterministic Policy Gradient

Let us continue to derive policy gradient with a **deterministic** policy $\mu$

$$\mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \Pr\big(s_{n+1} \mid s_0; \theta_\mu\big)$$

$$\Pr\big(s_{n+1} \mid s_0; \theta_\mu\big) = \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

# Deterministic Policy Gradient

Let us continue to derive policy gradient with a **deterministic** policy $\mu$

$$\mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \Pr\big(s_{n+1} \mid s_0; \theta_\mu\big)$$

$$\Pr\big(s_{n+1} \mid s_0; \theta_\mu\big) = \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

Plug state distribution into expected return

# Deterministic Policy Gradient

Let us continue to derive policy gradient with a **deterministic** policy $\mu$

$$\mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \Pr\big(s_{n+1} \mid s_0; \theta_\mu\big)$$

$$\Pr\big(s_{n+1} \mid s_0; \theta_\mu\big) = \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

Plug state distribution into expected return

$$\mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

# Deterministic Policy Gradient

$$\mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

# Deterministic Policy Gradient

$$\mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \dots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

Take gradient of both sides

# Deterministic Policy Gradient

$$\mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] = \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1,\dots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

Take gradient of both sides

$$\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$$

$$= \nabla_{\theta_\mu} \left[ \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1,\dots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right]$$

# Deterministic Policy Gradient

$$\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$$

$$= \nabla_{\theta_\mu}\left[\sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1,\dots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)\right]$$

# Deterministic Policy Gradient

$$\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$$

$$= \nabla_{\theta_\mu} \left[ \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right]$$

Gradient of sum is sum of gradient, move the gradient inside the sums

# Deterministic Policy Gradient

$$\nabla_{\theta_\mu} \mathbb{E}\left[\mathcal{G}(\tau) \mid s_0; \theta_\mu\right]$$

$$= \nabla_{\theta_\mu}\left[\sum_{n=0}^{\infty}\gamma^n \sum_{s_{n+1}\in S}\mathcal{R}(s_{n+1}) \sum_{s_1,\ldots,s_n\in S}\prod_{t=0}^{n}\mathrm{Tr}\left(s_{t+1}\mid s_t, \mu(s_t, \theta_\mu)\right)\right]$$

Gradient of sum is sum of gradient, move the gradient inside the sums

$$= \sum_{n=0}^{\infty}\gamma^n \sum_{s_{n+1}\in S}\mathcal{R}(s_{n+1}) \sum_{s_1,\ldots,s_n\in S}\nabla_{\theta_\mu}\left[\prod_{t=0}^{n}\mathrm{Tr}\left(s_{t+1}\mid s_t, \mu(s_t, \theta_\mu)\right)\right]$$

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \ldots, s_n \in S} \nabla_{\theta_\mu} \left[ \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right]$$

# Deterministic Policy Gradient

$$
= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \ldots, s_n \in S} \nabla_{\theta_\mu} \left[ \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right]
$$

Recall the log trick

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \dots, s_n \in S} \nabla_{\theta_\mu} \left[ \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right]$$

Recall the log trick

$$\nabla_x f(x) = f(x) \nabla_x \log f(x)$$

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \dots, s_n \in S} \nabla_{\theta_\mu} \left[ \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right]$$

Recall the log trick $\qquad\qquad \nabla_x f(x) = f(x) \nabla_x \log f(x)$

Apply log trick to product terms

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \ldots, s_n \in S} \nabla_{\theta_\mu} \left[ \prod_{t=0}^{n} \mathrm{Tr}\left(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\right) \right]$$

Recall the log trick $\qquad \nabla_x f(x) = f(x) \nabla_x \log f(x)$

Apply log trick to product terms

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \ldots, s_n \in S}$$

$$\prod_{t=0}^{n} \mathrm{Tr}\left(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\right) \nabla_{\theta_\mu} \left[ \log \left( \prod_{t=0}^{n} \mathrm{Tr}\left(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\right) \right) \right]$$

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1,\ldots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\nabla_{\theta_\mu} \left[ \log \left( \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right) \right]$$

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \dots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\nabla_{\theta_\mu} \left[ \log \left( \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right) \right]$$

Log of products is sum of logs

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\nabla_{\theta_\mu} \left[ \log \left( \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right) \right]$$

Log of products is sum of logs

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\nabla_{\theta_\mu} \left[ \sum_{t=0}^{n} \log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right]$$

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1,\ldots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\nabla_{\theta_\mu} \left[ \sum_{t=0}^{n} \log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right]$$

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1,\ldots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\nabla_{\theta_\mu} \left[ \sum_{t=0}^{n} \log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right]$$

Move the gradient inside sum

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1,\ldots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\nabla_{\theta_\mu} \left[ \sum_{t=0}^{n} \log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big) \right]$$

Move the gradient inside sum

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1,\ldots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\sum_{t=0}^{n} \nabla_{\theta_\mu} \log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\left(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\right)$$

$$\sum_{t=0}^{n} \nabla_{\theta_\mu} \log \mathrm{Tr}\left(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\right)$$

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1,\ldots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\sum_{t=0}^{n} \nabla_{\theta_\mu} \log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

Use the chain rule $\qquad\qquad\qquad\qquad \nabla_x f(g(x)) = \nabla_g[f(g(x))]\nabla_x g(x)$

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\sum_{t=0}^{n} \nabla_{\theta_\mu} \log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

Use the chain rule
$$\nabla_x f(g(x)) = \nabla_g[f(g(x))] \nabla_x g(x)$$

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\sum_{t=0}^{n} \nabla_\mu \big[\log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)\big] \cdot \nabla_{\theta_\mu} \mu(s_t, \theta_\mu)$$

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\sum_{t=0}^{n} \nabla_\mu \big[\log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)\big] \cdot \nabla_{\theta_\mu} \mu(s_t, \theta_\mu)$$

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \sum_{s_1, \dots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\sum_{t=0}^{n} \nabla_\mu \big[\log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)\big] \cdot \nabla_{\theta_\mu} \mu(s_t, \theta_\mu)$$

**Question:** Any problems?

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \sum_{s_1,\ldots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\sum_{t=0}^{n} \nabla_\mu \big[\log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)\big] \cdot \nabla_{\theta_\mu} \mu(s_t, \theta_\mu)$$

**Question:** Any problems?

We must know $\nabla \mathrm{Tr}$ to find the deterministic policy gradient

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \sum_{s_1, \ldots, s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\sum_{t=0}^{n} \nabla_\mu \big[\log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)\big] \cdot \nabla_{\theta_\mu} \mu(s_t, \theta_\mu)$$

**Question:** Any problems?

We must know $\nabla \mathrm{Tr}$ to find the deterministic policy gradient

In model-free RL, we do not know $\nabla \mathrm{Tr}$ – it might not be differentiable!

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \sum_{s_1,\ldots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\sum_{t=0}^{n} \nabla_\mu \big[\log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)\big] \cdot \nabla_{\theta_\mu} \mu(s_t, \theta_\mu)$$

**Question:** Any problems?

We must know $\nabla \mathrm{Tr}$ to find the deterministic policy gradient

In model-free RL, we do not know $\nabla \mathrm{Tr}$ – it might not be differentiable!

Stochastic policy gradient is very special – we do not need $\nabla \mathrm{Tr}$

# Deterministic Policy Gradient

$$= \sum_{n=0}^{\infty} \gamma^n \sum_{s_{n+1} \in S} \mathcal{R}(s_{n+1}) \cdot \sum_{s_1,\ldots,s_n \in S} \prod_{t=0}^{n} \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

$$\sum_{t=0}^{n} \nabla_\mu \big[\log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)\big] \cdot \nabla_{\theta_\mu} \mu(s_t, \theta_\mu)$$

**Question:** Any problems?

We must know $\nabla \mathrm{Tr}$ to find the deterministic policy gradient

In model-free RL, we do not know $\nabla \mathrm{Tr}$ – it might not be differentiable!

Stochastic policy gradient is very special – we do not need $\nabla \mathrm{Tr}$

Let me explain what I mean

# Deterministic Policy Gradient

With deterministic policy, $\mu$ inside Tr means chain rule

# Deterministic Policy Gradient

With deterministic policy, $\mu$ inside Tr means chain rule

$$\sum_{t=0}^{n} \nabla_{\theta_\mu} \log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

# Deterministic Policy Gradient

With deterministic policy, $\mu$ inside Tr means chain rule

$$\sum_{t=0}^{n} \nabla_{\theta_{\mu}} \log \mathrm{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_{\mu})\big)$$

With stochastic policy, $\pi$ **not inside** Tr (product rule)

# Deterministic Policy Gradient

With deterministic policy, $\mu$ inside Tr means chain rule

$$\sum_{t=0}^{n} \nabla_{\theta_\mu} \log \text{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

With stochastic policy, $\pi$ **not inside** Tr (product rule)

$$\sum_{t=0}^{n} \sum_{a_t \in A} \nabla_{\theta_\pi} \big[\text{Tr}(s_{t+1} \mid s_t, a_t) \cdot \log \pi(a_t \mid s_t; \theta_\pi)\big]$$

# Deterministic Policy Gradient

With deterministic policy, $\mu$ inside Tr means chain rule

$$\sum_{t=0}^{n} \nabla_{\theta_\mu} \log \text{Tr}\big(s_{t+1} \mid s_t, \mu(s_t, \theta_\mu)\big)$$

With stochastic policy, $\pi$ **not inside** Tr (product rule)

$$\sum_{t=0}^{n} \sum_{a_t \in A} \nabla_{\theta_\pi} \big[\text{Tr}\big(s_{t+1} \mid s_t, a_t\big) \cdot \log \pi(a_t \mid s_t; \theta_\pi)\big]$$

We use stochastic policies in RL because of this

# Deep Deterministic Policy Gradient

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$$

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] \qquad \theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$$

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] \qquad \theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$$

We failed – a deterministic policy gradient does not seem to work

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] \qquad \theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$$

We failed – a deterministic policy gradient does not seem to work

Can we try and optimize something else?

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] \qquad \theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$$

We failed – a deterministic policy gradient does not seem to work

Can we try and optimize something else?

**Question:** Could we replace $\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$ with something else?

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] \qquad \theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$$

We failed – a deterministic policy gradient does not seem to work

Can we try and optimize something else?

**Question:** Could we replace $\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$ with something else?

$$\nabla_{\theta_\mu} V\big(s_0, \theta_\mu\big)$$

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_{\mu}} \mathbb{E}\big[\mathcal{G}(\tau) \mid s_0; \theta_{\mu}\big] \qquad \theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_{\mu}} \mathbb{E}\big[\mathcal{G}(\tau) \mid s_0; \theta_{\mu}\big]$$

We failed – a deterministic policy gradient does not seem to work

Can we try and optimize something else?

**Question:** Could we replace $\nabla_{\theta_{\mu}} \mathbb{E}\big[\mathcal{G}(\tau) \mid s_0; \theta_{\mu}\big]$ with something else?

$$\nabla_{\theta_{\mu}} V\big(s_0, \theta_{\mu}\big) \qquad\qquad \nabla_{\theta_{\mu}} Q\big(s_0, a_0, \theta_{\mu}\big)$$

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] \qquad \theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$$

We failed – a deterministic policy gradient does not seem to work

Can we try and optimize something else?

**Question:** Could we replace $\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$ with something else?

$$\nabla_{\theta_\mu} V\big(s_0, \theta_\mu\big) \qquad\qquad \nabla_{\theta_\mu} Q\big(s_0, a_0, \theta_\mu\big)$$

Let us try to derive deterministic policy gradient again

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big] \qquad \theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$$

We failed – a deterministic policy gradient does not seem to work

Can we try and optimize something else?

**Question:** Could we replace $\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$ with something else?

$$\nabla_{\theta_\mu} V\big(s_0, \theta_\mu\big) \qquad\qquad \nabla_{\theta_\mu} Q\big(s_0, a_0, \theta_\mu\big)$$

Let us try to derive deterministic policy gradient again

This time, take gradient of $V$ instead of gradient of $\mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]$

# Deep Deterministic Policy Gradient

$$V\left(s_0, \theta_\mu\right) = \mathbb{E}\left[\mathcal{R}(s_1) \mid s_0; \theta_\mu\right] + \gamma V\left(s_0, \theta_\mu\right)$$

# Deep Deterministic Policy Gradient

$$V(s_0, \theta_\mu) = \mathbb{E}\big[\mathcal{R}(s_1) \mid s_0; \theta_\mu\big] + \gamma V(s_0, \theta_\mu)$$

Rewrite $V$ in terms of $Q$

# Deep Deterministic Policy Gradient

$$V\big(s_0, \theta_\mu\big) = \mathbb{E}\big[\mathcal{R}(s_1) \mid s_0; \theta_\mu\big] + \gamma V\big(s_0, \theta_\mu\big)$$

Rewrite $V$ in terms of $Q$

$$V\big(s_0, \theta_\mu\big) = Q\big(s_0, a, \theta_\mu\big); \quad a = \mu\big(s_1, \theta_\mu\big)$$

# Deep Deterministic Policy Gradient

$$V(s_0, \theta_\mu) = \mathbb{E}\big[\mathcal{R}(s_1) \mid s_0; \theta_\mu\big] + \gamma V(s_0, \theta_\mu)$$

Rewrite $V$ in terms of $Q$

$$V(s_0, \theta_\mu) = Q(s_0, a, \theta_\mu); \quad a = \mu(s_1, \theta_\mu)$$

Plug in $\mu$ for $a$

# Deep Deterministic Policy Gradient

$$V(s_0, \theta_\mu) = \mathbb{E}\big[\mathcal{R}(s_1) \mid s_0; \theta_\mu\big] + \gamma V(s_0, \theta_\mu)$$

Rewrite $V$ in terms of $Q$

$$V(s_0, \theta_\mu) = Q(s_0, a, \theta_\mu); \quad a = \mu(s_1, \theta_\mu)$$

Plug in $\mu$ for $a$

$$V(s_0, \theta_\mu) = Q(s_0, \mu(s_0, \theta_\mu), \theta_\mu)$$

# Deep Deterministic Policy Gradient

$$V(s_0, \theta_\mu) = \mathbb{E}\big[\mathcal{R}(s_1) \mid s_0; \theta_\mu\big] + \gamma V(s_0, \theta_\mu)$$

Rewrite $V$ in terms of $Q$

$$V(s_0, \theta_\mu) = Q(s_0, a, \theta_\mu); \quad a = \mu(s_1, \theta_\mu)$$

Plug in $\mu$ for $a$

$$V(s_0, \theta_\mu) = Q(s_0, \mu(s_0, \theta_\mu), \theta_\mu)$$

Take the gradient of both sides

# Deep Deterministic Policy Gradient

$$V(s_0, \theta_\mu) = \mathbb{E}\big[\mathcal{R}(s_1) \mid s_0; \theta_\mu\big] + \gamma V(s_0, \theta_\mu)$$

Rewrite $V$ in terms of $Q$

$$V(s_0, \theta_\mu) = Q(s_0, a, \theta_\mu); \quad a = \mu(s_1, \theta_\mu)$$

Plug in $\mu$ for $a$

$$V(s_0, \theta_\mu) = Q(s_0, \mu(s_0, \theta_\mu), \theta_\mu)$$

Take the gradient of both sides

$$\nabla_{\theta_\mu} V(s_0, a_0, \theta_\mu) = \nabla_{\theta_\mu} Q(s_0, \mu(s_0, \theta_\mu), \theta_\mu)$$

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_{\theta_\mu} Q(s_0, \mu(s_0, \theta_\mu), \theta_\mu)$$

Use the chain rule
$$\nabla_x f(g(x)) = \nabla_g [f(g(x))] \nabla_x g(x)$$

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_{\theta_\mu} Q(s_0, \mu(s_0, \theta_\mu), \theta_\mu)$$

Use the chain rule $\qquad\qquad \nabla_x f(g(x)) = \nabla_g [f(g(x))] \nabla_x g(x)$

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_\mu Q(s_0, \mu(s_0, \theta_\mu), \theta_\mu) \cdot \nabla_{\theta_\mu} \mu(s_0, \theta_\mu)$$

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_\mu Q(s_1, \mu(s_1, \theta_\mu), \theta_\mu) \cdot \nabla_{\theta_\mu} \mu(s_1, \theta_\mu)$$
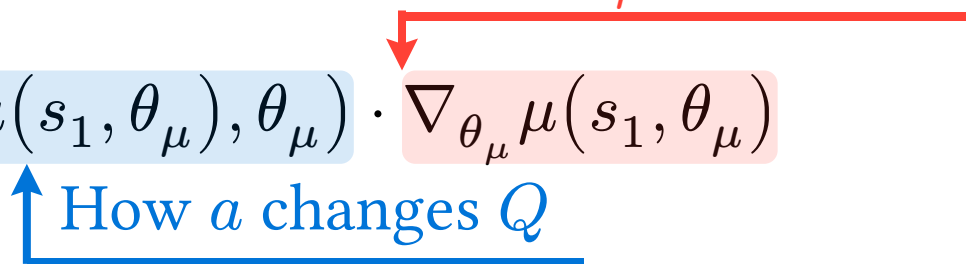
# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_\mu Q(s_1, \mu(s_1, \theta_\mu), \theta_\mu) \cdot \nabla_{\theta_\mu} \mu(s_1, \theta_\mu)$$

Let us inspect these terms more closely

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_\mu Q(s_1, \mu(s_1, \theta_\mu), \theta_\mu) \cdot \nabla_{\theta_\mu} \mu(s_1, \theta_\mu)$$

Let us inspect these terms more closely

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_\mu Q\big(s_1, \mu(s_1, \theta_\mu), \theta_\mu\big) \cdot \nabla_{\theta_\mu} \mu(s_1, \theta_\mu)$$

How $a$ changes $Q$

Let us inspect these terms more closely

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_\mu Q(s_1, \mu(s_1, \theta_\mu), \theta_\mu) \cdot \nabla_{\theta_\mu} \mu(s_1, \theta_\mu)$$

How $a$ changes $Q$

Let us inspect these terms more closely

The second term is easy to compute, gradient of deterministic policy

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_\mu Q(s_1, \mu(s_1, \theta_\mu), \theta_\mu) \cdot \nabla_{\theta_\mu} \mu(s_1, \theta_\mu)$$
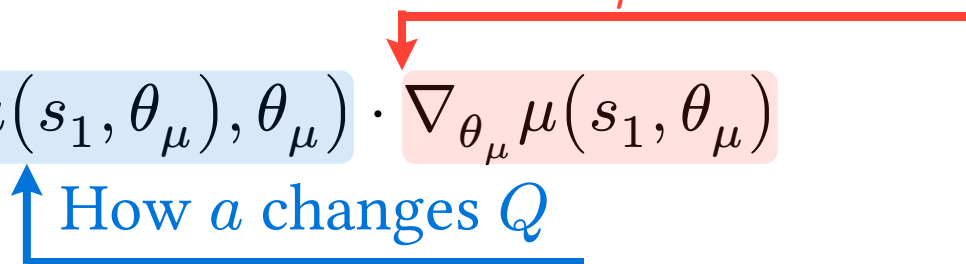
How $a$ changes $Q$

Let us inspect these terms more closely

The second term is easy to compute, gradient of deterministic policy

Analytical solution for the first term is difficult (recursive definition)

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_\mu Q(s_1, \mu(s_1, \theta_\mu), \theta_\mu) \cdot \nabla_{\theta_\mu} \mu(s_1, \theta_\mu)$$

How $\theta_\mu$ changes $a$

How $a$ changes $Q$

Let us inspect these terms more closely

The second term is easy to compute, gradient of deterministic policy

Analytical solution for the first term is difficult (recursive definition)
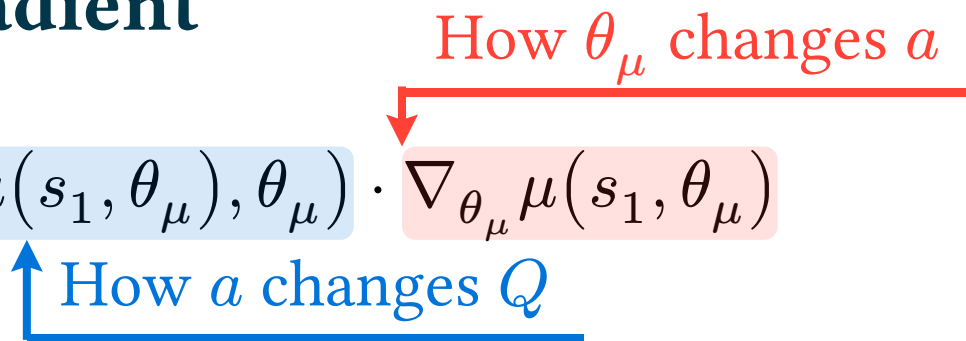
But what if $Q$ is a neural network?

# Deep Deterministic Policy Gradient

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_\mu Q(s_1, \mu(s_1, \theta_\mu), \theta_\mu) \cdot \nabla_{\theta_\mu} \mu(s_1, \theta_\mu)$$

How $a$ changes $Q$

Let us inspect these terms more closely

The second term is easy to compute, gradient of deterministic policy

Analytical solution for the first term is difficult (recursive definition)

But what if $Q$ is a neural network?

We can backpropagate through $Q$ without worrying about recursion

# Deep Deterministic Policy Gradient

How $\theta_\mu$ changes $a$

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_\mu Q(s_1, \mu(s_1, \theta_\mu), \theta_\mu) \cdot \nabla_{\theta_\mu} \mu(s_1, \theta_\mu)$$

How $a$ changes $Q$

# Deep Deterministic Policy Gradient

How $\theta_\mu$ changes $a$

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \underbrace{\nabla_\mu Q(s_1, \mu(s_1, \theta_\mu), \theta_\mu)}_{\text{How } a \text{ changes } Q} \cdot \nabla_{\theta_\mu} \mu(s_1, \theta_\mu)$$

Writing the code makes it look easy

# Deep Deterministic Policy Gradient

How $\theta_\mu$ changes $a$

$$\nabla_{\theta_\mu} V(s_0, \theta_\mu) = \nabla_\mu Q(s_1, \mu(s_1, \theta_\mu), \theta_\mu) \cdot \nabla_{\theta_\mu} \mu(s_1, \theta_\mu)$$

How $a$ changes $Q$

Writing the code makes it look easy

```
def V(s, Q_nn, mu_nn):
    a = mu_nn(s)
    return Q_nn(s, a)


# Learn the policy that maximizes V
# Make sure to differentiate w.r.t policy parameters!
J = grad(V, argnums=2)(states, Q_nn, mu_nn)
mu_nn = optimizer.update(mu_nn, J)
```

# Deep Deterministic Policy Gradient

This assumes we know the $Q$ function for $\theta_\mu$

# Deep Deterministic Policy Gradient

This assumes we know the $Q$ function for $\theta_\mu$

We can learn $Q$ for any policy (before we focused on greedy)

# Deep Deterministic Policy Gradient

This assumes we know the $Q$ function for $\theta_\mu$

We can learn $Q$ for any policy (before we focused on greedy)

$$Q(s_0, a_0, \theta_\mu, \theta_Q) = \mathbb{E}[\mathcal{R}(s_1) \mid s_0, a_0] + \gamma Q(s_1, \mu(s_1, \theta_\mu), \theta_\mu, \theta_Q)$$

# Deep Deterministic Policy Gradient

This assumes we know the $Q$ function for $\theta_\mu$

We can learn $Q$ for any policy (before we focused on greedy)

$$Q(s_0, a_0, \theta_\mu, \theta_Q) = \mathbb{E}[\mathcal{R}(s_1) \mid s_0, a_0] + \gamma Q(s_1, \mu(s_1, \theta_\mu), \theta_\mu, \theta_Q)$$

```python
def V(s, Q_nn, mu_nn):
    a = mu_nn(s)
    return Q_nn(s, a)
# Before, we learned policy params to maximize Q
# Now, we learn params of Q following policy (argnums=2)
J = grad(V, argnums=1)(states, Q_nn, mu_nn)
Q_nn = optimizer.update(Q_nn, J)
```

# Deep Deterministic Policy Gradient

**Definition:** Deep Deterministic Policy Gradient (DDPG) decomposes $V$ into a deterministic policy $\mu$ and $Q$, learning them jointly

# Deep Deterministic Policy Gradient

**Definition:** Deep Deterministic Policy Gradient (DDPG) decomposes $V$ into a deterministic policy $\mu$ and $Q$, learning them jointly

**Step 1:** Learn a $Q$ function for $\theta_\mu$

# Deep Deterministic Policy Gradient

**Definition:** Deep Deterministic Policy Gradient (DDPG) decomposes $V$ into a deterministic policy $\mu$ and $Q$, learning them jointly

**Step 1:** Learn a $Q$ function for $\theta_\mu$

$$\theta_{Q,i+1} = \arg\min_{\theta_{Q,i}}$$

$$\left( Q(s_0, a_0, \theta_{\mu,i}, \theta_{Q,i}) - \left( \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + \gamma Q(s_1, \mu(s_1, \theta_{\mu,i}), \theta_{\mu,i}, \theta_{Q,i}) \right) \right)$$

# Deep Deterministic Policy Gradient

**Definition:** Deep Deterministic Policy Gradient (DDPG) decomposes $V$ into a deterministic policy $\mu$ and $Q$, learning them jointly

**Step 1:** Learn a $Q$ function for $\theta_\mu$

$$\theta_{Q,i+1} = \arg \min_{\theta_{Q,i}}$$

$$\Big( Q(s_0, a_0, \theta_{\mu,i}, \theta_{Q,i}) - \Big( \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + \gamma Q(s_1, \mu(s_1, \theta_{\mu,i}), \theta_{\mu,i}, \theta_{Q,i}) \Big) \Big)$$

**Step 2:** Learn $\theta_\mu$ that maximizes $Q$

# Deep Deterministic Policy Gradient

**Definition:** Deep Deterministic Policy Gradient (DDPG) decomposes $V$ into a deterministic policy $\mu$ and $Q$, learning them jointly

**Step 1:** Learn a $Q$ function for $\theta_\mu$

$$\theta_{Q,i+1} = \arg\min_{\theta_{Q,i}}$$

$$\left( Q(s_0, a_0, \theta_{\mu,i}, \theta_{Q,i}) - \left( \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + \gamma Q(s_1, \mu(s_1, \theta_{\mu,i}), \theta_{\mu,i}, \theta_{Q,i}) \right) \right)$$

**Step 2:** Learn $\theta_\mu$ that maximizes $Q$

$$\theta_{\mu,i+1} = \theta_{\mu,i} + \alpha \cdot Q\left( s_0, \mu(s_0, \theta_\mu), \theta_{\mu,i}, \theta_{Q,i+1} \right)$$

# Deep Deterministic Policy Gradient

**Definition:** Deep Deterministic Policy Gradient (DDPG) decomposes $V$ into a deterministic policy $\mu$ and $Q$, learning them jointly

**Step 1:** Learn a $Q$ function for $\theta_\mu$

$$\theta_{Q,i+1} = \arg \min_{\theta_{Q,i}}$$

$$\left( Q(s_0, a_0, \theta_{\mu,i}, \theta_{Q,i}) - \left( \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + \gamma Q(s_1, \mu(s_1, \theta_{\mu,i}), \theta_{\mu,i}, \theta_{Q,i}) \right) \right)$$

**Step 2:** Learn $\theta_\mu$ that maximizes $Q$

$$\theta_{\mu,i+1} = \theta_{\mu,i} + \alpha \cdot Q\left( s_0, \mu(s_0, \theta_\mu), \theta_{\mu,i}, \theta_{Q,i+1} \right)$$

Repeat until convergence, $\theta_{\mu,i+1} = \theta_{\mu,i}, \quad \theta_{Q,i+1} = \theta_{Q,i}$

# Deep Deterministic Policy Gradient

$$\theta_{Q,i+1} = \arg \min_{\theta_{Q,i}}$$

$$\left( Q(s_0, a_0, \theta_{\mu,i}, \theta_{Q,i}) - \left( \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + \gamma Q(s_1, \mu(s_1, \theta_{\mu,i}), \theta_{\mu,i}, \theta_{Q,i}) \right) \right)$$

$$\theta_{\mu,i+1} = \theta_{\mu,i} + \alpha \cdot Q(s_0, \mu(s_0, \theta_\mu), \theta_{\mu,i}, \theta_{Q,i+1})$$

# Deep Deterministic Policy Gradient

$$\theta_{Q,i+1} = \arg\min_{\theta_{Q,i}}$$

$$\Big( Q(s_0, a_0, \theta_{\mu,i}, \theta_{Q,i}) - \big( \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + \gamma Q(s_1, \mu(s_1, \theta_{\mu,i}), \theta_{\mu,i}, \theta_{Q,i}) \big) \Big)$$

$$\theta_{\mu,i+1} = \theta_{\mu,i} + \alpha \cdot Q(s_0, \mu(s_0, \theta_\mu), \theta_{\mu,i}, \theta_{Q,i+1})$$

**Question:** Is DDPG on-policy or off-policy?

# Deep Deterministic Policy Gradient

$$\theta_{Q,i+1} = \arg\min_{\theta_{Q,i}}$$

$$\Big(Q(s_0, a_0, \theta_{\mu,i}, \theta_{Q,i}) - \big(\hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + \gamma Q(s_1, \mu(s_1, \theta_{\mu,i}), \theta_{\mu,i}, \theta_{Q,i})\big)\Big)$$

$$\theta_{\mu,i+1} = \theta_{\mu,i} + \alpha \cdot Q\big(s_0, \mu(s_0, \theta_\mu), \theta_{\mu,i}, \theta_{Q,i+1}\big)$$

**Question:** Is DDPG on-policy or off-policy?

**Answer:** Depends on how we train $Q$:

# Deep Deterministic Policy Gradient

$$\theta_{Q,i+1} = \underset{\theta_{Q,i}}{\arg\min}$$

$$\Big( Q(s_0, a_0, \theta_{\mu,i}, \theta_{Q,i}) - \Big( \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + \gamma Q(s_1, \mu(s_1, \theta_{\mu,i}), \theta_{\mu,i}, \theta_{Q,i}) \Big) \Big)$$

$$\theta_{\mu,i+1} = \theta_{\mu,i} + \alpha \cdot Q(s_0, \mu(s_0, \theta_\mu), \theta_{\mu,i}, \theta_{Q,i+1})$$

**Question:** Is DDPG on-policy or off-policy?

**Answer:** Depends on how we train $Q$:

- TD $Q$ is off-policy

# Deep Deterministic Policy Gradient

$$\theta_{Q,i+1} = \arg\min_{\theta_{Q,i}}$$

$$\Big( Q(s_0, a_0, \theta_{\mu,i}, \theta_{Q,i}) - \big( \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + \gamma Q(s_1, \mu(s_1, \theta_{\mu,i}), \theta_{\mu,i}, \theta_{Q,i}) \big) \Big)$$

$$\theta_{\mu,i+1} = \theta_{\mu,i} + \alpha \cdot Q\big(s_0, \mu(s_0, \theta_\mu), \theta_{\mu,i}, \theta_{Q,i+1}\big)$$

**Question:** Is DDPG on-policy or off-policy?

**Answer:** Depends on how we train $Q$:

- TD $Q$ is off-policy
- MC $Q$ is on-policy

# Deep Deterministic Policy Gradient

$$\theta_{Q,i+1} = \arg \min_{\theta_{Q,i}}$$

$$\Big( Q(s_0, a_0, \theta_{\mu,i}, \theta_{Q,i}) - \big( \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + \gamma Q(s_1, \mu(s_1, \theta_{\mu,i}), \theta_{\mu,i}, \theta_{Q,i}) \big) \Big)$$

$$\theta_{\mu,i+1} = \theta_{\mu,i} + \alpha \cdot Q\big( s_0, \mu(s_0, \theta_\mu), \theta_{\mu,i}, \theta_{Q,i+1} \big)$$

**Question:** Is DDPG on-policy or off-policy?

**Answer:** Depends on how we train $Q$:

- TD $Q$ is off-policy
- MC $Q$ is on-policy

Almost **all** good off-policy actor-critic algorithms are based on DDPG

# Deep Deterministic Policy Gradient

**Summary:**

# Deep Deterministic Policy Gradient

**Summary:**

- Wanted to extend Q learning to continuous action spaces

# Deep Deterministic Policy Gradient

**Summary:**

- Wanted to extend Q learning to continuous action spaces
  - Simple greedy policy does not work!

# Deep Deterministic Policy Gradient

**Summary:**

- Wanted to extend Q learning to continuous action spaces
  - ▸ Simple greedy policy does not work!
- Introduced learned deterministic policy $a = \mu(s, \theta_\mu)$

# Deep Deterministic Policy Gradient

## Summary:

- Wanted to extend Q learning to continuous action spaces
    - ‣ Simple greedy policy does not work!
- Introduced learned deterministic policy $a = \mu(s, \theta_\mu)$
- Failed to find deterministic policy gradient $\nabla_{\theta_\mu} \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0 ; \theta_\mu\big]$

# Deep Deterministic Policy Gradient

**Summary:**

- Wanted to extend Q learning to continuous action spaces
  - ▸ Simple greedy policy does not work!
- Introduced learned deterministic policy $a = \mu\left(s, \theta_\mu\right)$
- Failed to find deterministic policy gradient $\nabla_{\theta_\mu} \mathbb{E}\left[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\right]$
  - ▸ Must know $\nabla \operatorname{Tr}$

# Deep Deterministic Policy Gradient

## Summary:

- Wanted to extend Q learning to continuous action spaces
  - ‣ Simple greedy policy does not work!
- Introduced learned deterministic policy $a = \mu(s, \theta_\mu)$
- Failed to find deterministic policy gradient $\nabla_{\theta_\mu} \mathbb{E}\left[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\right]$
  - ‣ Must know $\nabla \operatorname{Tr}$
- Try gradient ascent on $\nabla_{\theta_\mu} V$ instead of $\nabla_{\theta_\mu} \mathbb{E}\left[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\right]$
  - ‣ Trick: Factor $V$ in terms of $Q, \mu$
  - ‣ Iterative optimization of $\theta_Q$ and $\theta_\mu$

# Deep Deterministic Policy Gradient

**Summary:**

- Wanted to extend Q learning to continuous action spaces
  - ‣ Simple greedy policy does not work!
- Introduced learned deterministic policy $a = \mu(s, \theta_\mu)$
- Failed to find deterministic policy gradient $\nabla_{\theta_\mu} \mathbb{E}\left[\mathcal{G}(\tau) \mid s_0; \theta_\mu\right]$
  - ‣ Must know $\nabla \text{Tr}$
- Try gradient ascent on $\nabla_{\theta_\mu} V$ instead of $\nabla_{\theta_\mu} \mathbb{E}\left[\mathcal{G}(\tau) \mid s_0; \theta_\mu\right]$
  - ‣ Trick: Factor $V$ in terms of $Q, \mu$
  - ‣ Iterative optimization of $\theta_Q$ and $\theta_\mu$

# Deep Deterministic Policy Gradient

Another way to think of DDPG:

# Deep Deterministic Policy Gradient

Another way to think of DDPG:

$$\max_{\theta_\mu} Q\big(s, \mu(s, \theta_\mu), \theta_\mu\big) \approx \max_{a \in A} Q\big(s, a, \theta_\mu\big)$$

# Deep Deterministic Policy Gradient

Another way to think of DDPG:

$$\max_{\theta_\mu} Q\big(s, \mu(s, \theta_\mu), \theta_\mu\big) \approx \max_{a \in A} Q\big(s, a, \theta_\mu\big)$$

$\mu$ is a continuous relaxation of the greedy policy

# Deep Deterministic Policy Gradient

Another way to think of DDPG:

$$\max_{\theta_\mu} Q\big(s, \mu(s, \theta_\mu), \theta_\mu\big) \approx \max_{a \in A} Q\big(s, a, \theta_\mu\big)$$

$\mu$ is a continuous relaxation of the greedy policy

Because $\mu$ is a neural network, it can generalize to continuous $s, a$

# Deep Deterministic Policy Gradient

One small problem: deterministic policy means no exploration!

Actor Critic II

# Deep Deterministic Policy Gradient

One small problem: deterministic policy means no exploration!

We must visit other states/actions to find optimal $\theta_Q, \theta_\mu$

# Deep Deterministic Policy Gradient

One small problem: deterministic policy means no exploration!

We must visit other states/actions to find optimal $\theta_Q, \theta_\mu$

With Q learning, we had epsilon greedy policy

# Deep Deterministic Policy Gradient

One small problem: deterministic policy means no exploration!

We must visit other states/actions to find optimal $\theta_Q, \theta_\mu$

With Q learning, we had epsilon greedy policy

$$\pi(a \mid s; \theta_\pi) = \begin{cases} 1 - \varepsilon : a = \arg\max_{a \in A} Q(s, a, \theta_\pi) \\ \varepsilon \qquad : a = \text{uniform}(A) \end{cases}$$

# Deep Deterministic Policy Gradient

One small problem: deterministic policy means no exploration!

We must visit other states/actions to find optimal $\theta_Q, \theta_\mu$

With Q learning, we had epsilon greedy policy

$$\pi(a \mid s; \theta_\pi) = \begin{cases} 1 - \varepsilon : a = \arg \max_{a \in A} Q(s, a, \theta_\pi) \\ \varepsilon \qquad : a = \text{uniform}(A) \end{cases}$$

**Question:** What about DDPG exploration? HINT: Continuous actions

# Deep Deterministic Policy Gradient

One small problem: deterministic policy means no exploration!

We must visit other states/actions to find optimal $\theta_Q, \theta_\mu$

With Q learning, we had epsilon greedy policy

$$\pi(a \mid s; \theta_\pi) = \begin{cases} 1 - \varepsilon : a = \arg\max_{a \in A} Q(s, a, \theta_\pi) \\ \varepsilon \qquad : a = \text{uniform}(A) \end{cases}$$

**Question:** What about DDPG exploration? HINT: Continuous actions

$$\pi(a \mid s; \theta_\mu) = \begin{cases} 1 - \varepsilon : a = Q\big(s, \mu(s, \theta_\mu), \theta_\mu\big) \\ \varepsilon \qquad : a = \text{uniform}(A) \end{cases}$$

# Deep Deterministic Policy Gradient

In practice, we can do something a little smarter

# Deep Deterministic Policy Gradient

In practice, we can do something a little smarter

Uniform actions cover the full action space, but can take a while to learn

# Deep Deterministic Policy Gradient

In practice, we can do something a little smarter

Uniform actions cover the full action space, but can take a while to learn

We can add noise to a good action instead of completely random actions

# Deep Deterministic Policy Gradient

In practice, we can do something a little smarter

Uniform actions cover the full action space, but can take a while to learn

We can add noise to a good action instead of completely random actions

$$\pi(a \mid s; \theta_\mu) = \mu(s, \theta_\mu) + \text{Normal}(0, \sigma)$$

# Deep Deterministic Policy Gradient

In practice, we can do something a little smarter

Uniform actions cover the full action space, but can take a while to learn

We can add noise to a good action instead of completely random actions

$$\pi\big(a \mid s; \theta_\mu\big) = \mu\big(s, \theta_\mu\big) + \text{Normal}(0, \sigma)$$

This tends to learn more quickly

# Deep Deterministic Policy Gradient

In practice, we can do something a little smarter

Uniform actions cover the full action space, but can take a while to learn

We can add noise to a good action instead of completely random actions

$$\pi(a \mid s; \theta_\mu) = \mu(s, \theta_\mu) + \text{Normal}(0, \sigma)$$

This tends to learn more quickly

Normal noise (infinite support) guarantees full action space coverage

# Coding

# Coding

Like policy gradient, the math and code is different

# Coding

Like policy gradient, the math and code is different

- Construct $\mu$ and $Q$ networks

# Coding

Like policy gradient, the math and code is different

- Construct $\mu$ and $Q$ networks
- Sample actions

# Coding

Like policy gradient, the math and code is different

- Construct $\mu$ and $Q$ networks
- Sample actions
  - ▸ Be careful that random actions are in action space!

# Coding

Like policy gradient, the math and code is different

- Construct $\mu$ and $Q$ networks
- Sample actions
  - ‣ Be careful that random actions are in action space!
  - ‣ $A = [0, 2\pi]$, then $a = 2.1\pi$ not ok

# Coding

Like policy gradient, the math and code is different

- Construct $\mu$ and $Q$ networks
- Sample actions
  - ▸ Be careful that random actions are in action space!
  - ▸ $A = [0, 2\pi]$, then $a = 2.1\pi$ not ok
- Iteratively train $\theta_Q, \theta_\mu$

# Coding

First, construct deterministic policy

# Coding

First, construct deterministic policy

```
mu = Sequential([
    Linear(state_size, hidden_size),
    Lambda(leaky_relu),
    Linear(hidden_size, hidden_size),
    Lambda(leaky_relu),
    Linear(hidden_size, action_dims),
])
```

# Coding

First, construct deterministic policy

```
mu = Sequential([
    Linear(state_size, hidden_size),
    Lambda(leaky_relu),
    Linear(hidden_size, hidden_size),
    Lambda(leaky_relu),
    Linear(hidden_size, action_dims),
])
```

Here, `action_dims` correspond to the number of continuous dimensions

# Coding

First, construct deterministic policy

```
mu = Sequential([
    Linear(state_size, hidden_size),
    Lambda(leaky_relu),
    Linear(hidden_size, hidden_size),
    Lambda(leaky_relu),
    Linear(hidden_size, action_dims),
])
```

Here, `action_dims` correspond to the number of continuous dimensions

BenBen: $A = [0, 2\pi]^{12}$, so `action_dims=12`

# Coding

Now, we need to make sure actions do not leave action space!

# Coding

Now, we need to make sure actions do not leave action space!

- BenBen: $A \in [0, 2\pi]^{12}$, `lower=[0, 0, ...]`, `upper=[2pi, 2pi, ...]`

# Coding

Now, we need to make sure actions do not leave action space!

- BenBen: $A \in [0, 2\pi]^{12}$, `lower=[0, 0, ...]`, `upper=[2pi, 2pi, ...]`

Squash actions to limits using tanh

# Coding

Now, we need to make sure actions do not leave action space!

- BenBen: $A \in [0, 2\pi]^{12}$, `lower=[0, 0, ...]`, `upper=[2pi, 2pi, ...]`

Squash actions to limits using tanh

- Can also use `clip(action, lower, upper)`, but this zeros gradients

# Coding

# Coding

Now, we need to make sure actions do not leave action space!

- BenBen: $A \in [0, 2\pi]^{12}$, `lower=[0, 0, ...]`, `upper=[2pi, 2pi, ...]`

Squash actions to limits using tanh

- Can also use `clip(action, lower, upper)`, but this zeros gradients
- Agent can get stuck taking limiting actions (e.g., $-0.01$ or $2.1\pi$)

# Coding

Now, we need to make sure actions do not leave action space!

- BenBen: $A \in [0, 2\pi]^{12}$, `lower=[0, 0, ...]`, `upper=[2pi, 2pi, ...]`

Squash actions to limits using tanh

- Can also use `clip(action, lower, upper)`, but this zeros gradients
- Agent can get stuck taking limiting actions (e.g., $-0.01$ or $2.1\pi$)

```python
def bound_action(action, lower, upper):
    return 0.5 * (upper + lower) + 0.5 * (upper - lower)
        * tanh(action)
```

# Coding

Now, we need to make sure actions do not leave action space!

- BenBen: $A \in [0, 2\pi]^{12}$, `lower=[0, 0, ...]`, `upper=[2pi, 2pi, ...]`

Squash actions to limits using tanh

- Can also use `clip(action, lower, upper)`, but this zeros gradients
- Agent can get stuck taking limiting actions (e.g., $-0.01$ or $2.1\pi$)

```python
def bound_action(action, lower, upper):
    return 0.5 * (upper + lower) + 0.5 * (upper - lower)
        * tanh(action)

def sample_action(mu, state, A_bounds, std):
    action = mu(state)
    noisy_action = action + normal(0, std) # Explore
    return bound_action(noisy_action, *A_bounds)
```

# Coding

Now construct $Q$

# Coding

Now construct $Q$

```
Q = Sequential([
    # Different from DQN network
    # Input action and state together
    Lambda(lambda s, a: concatenate(s, a)),
    Linear(state_size + action_dims, hidden_size),
    Lambda(leaky_relu),
    Linear(hidden_size, hidden_size),
    Lambda(leaky_relu),
    Linear(hidden_size, 1), # Single value for Q(s, a)
])
```

# Coding

```
while not terminated:
    # Exploration: make sure actions within action space!
    action = sample_action(mu, state, bounds, std)
    transition = env.step(action)
    replay_buffer.append(transition)
    data = replay_buffer.sample()
    # Theta_pi params are in mu neural network
    # Argnums tells us differentiation variable
    J_Q = grad(Q_loss, argnums=0)(theta_Q, theta_T, mu, data)
    theta_Q = apply_updates(J_Q, ...)
    J_mu = grad(mu_loss, argnums=0)(mu, theta_Q, data)
    mu = apply_updates(J_mu, ...)
    if step % 200 == 0: # Target network necessary
        theta_T = theta_Q
```

# Coding

```python
def Q_loss(theta_Q, theta_T, theta_pi, data):
    Qnet = combine(Q, theta_Q)
    Tnet = combine(Q, theta_T) # Target network
    # Predict Q values for action we took
    prediction = vmap(Qnet)(data.state, data.action)
    # Now compute labels using TD error
    next_action = vmap(mu)(data.next_state)
    # NOTE: No max! Mu approximates argmax
    next_Q = vmap(Tnet)(data.next_state, next_action)
    label = data.reward + gamma * data.done * next_Q
    return (prediction - label) ** 2
```

# Coding

```python
def mu_loss(mu, theta_Q, data):
    # Find the action that maximizes the Q function
    Qnet = combine(Q, theta_Q)
    # Instead of multiply, chain rule -- plug action into Q
    action = vmap(mu)(data.state)
    q_value = vmap(Qnet)(data.state, action)
    # Update the policy parameters to maximize the Q value
    # Gradient ascent but we min loss, use negative
    return -q_value
```

# Max Entropy RL

# Max Entropy RL

Many algorithms add improvements to DDPG

# Max Entropy RL

Many algorithms add improvements to DDPG

The most popular algorithm based on DDPG is **Soft Actor Critic** (SAC)

# Max Entropy RL

Many algorithms add improvements to DDPG

The most popular algorithm based on DDPG is **Soft Actor Critic** (SAC)

SAC is arguably the "best" model-free algorithm

# Max Entropy RL

Many algorithms add improvements to DDPG

The most popular algorithm based on DDPG is **Soft Actor Critic** (SAC)

SAC is arguably the "best" model-free algorithm

The motivation for SAC comes from **max-entropy RL**

# Max Entropy RL

Many algorithms add improvements to DDPG

The most popular algorithm based on DDPG is **Soft Actor Critic** (SAC)

SAC is arguably the "best" model-free algorithm

The motivation for SAC comes from **max-entropy RL**

We will very briefly cover max-entropy RL

# Max Entropy RL

Many algorithms add improvements to DDPG

The most popular algorithm based on DDPG is **Soft Actor Critic** (SAC)

SAC is arguably the "best" model-free algorithm

The motivation for SAC comes from **max-entropy RL**

We will very briefly cover max-entropy RL

First, let us introduce entropy

# Max Entropy RL

Entropy measures the uncertainty of a distribution

# Max Entropy RL
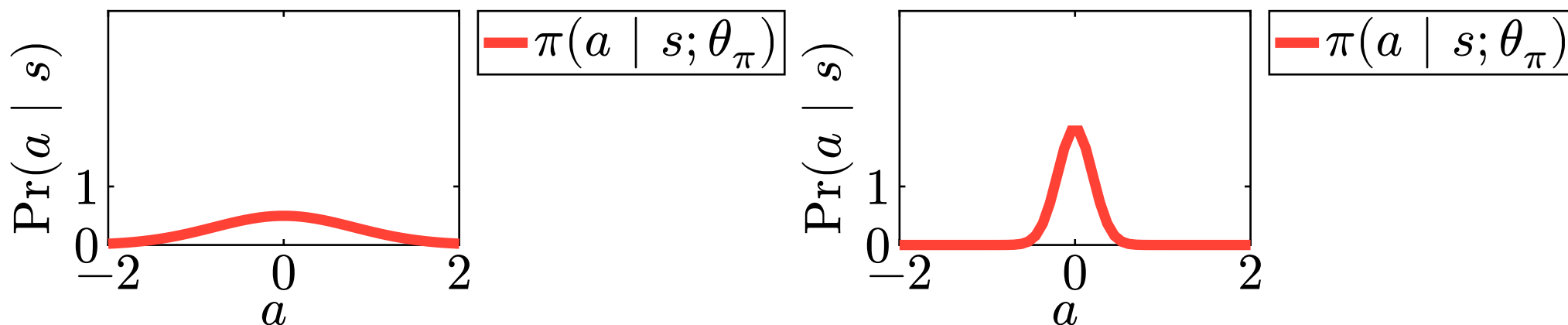
Entropy measures the uncertainty of a distribution

$$H(\pi(a \mid s; \theta_\pi))$$

# Max Entropy RL

Entropy measures the uncertainty of a distribution
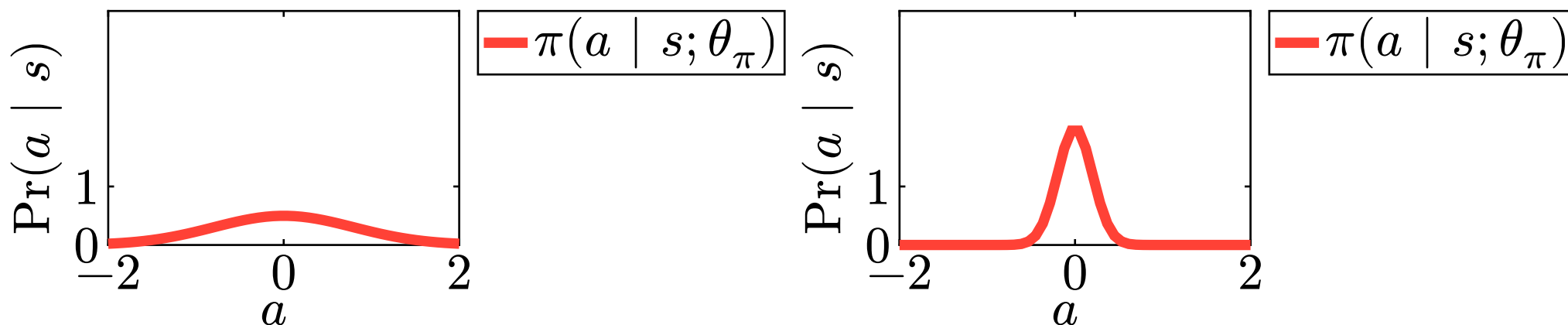
$$H(\pi(a \mid s; \theta_\pi))$$

# Max Entropy RL

Entropy measures the uncertainty of a distribution

$$H(\pi(a \mid s; \theta_\pi))$$

# Max Entropy RL

Entropy measures the uncertainty of a distribution
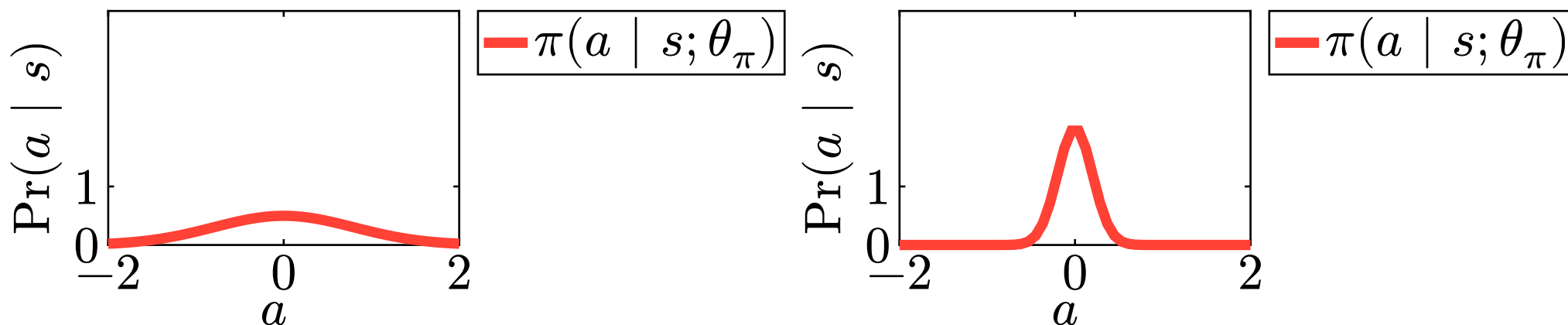
$$H(\pi(a \mid s; \theta_\pi))$$



**Question:** Which policy has higher entropy?

# Max Entropy RL

Entropy measures the uncertainty of a distribution

$$H(\pi(a \mid s; \theta_\pi))$$



**Question:** Which policy has higher entropy?

Left policy, more uncertain/random

# Max Entropy RL

In maximum entropy RL, we change the objective

# Max Entropy RL

In maximum entropy RL, we change the objective

$$\arg\max_{\theta_\pi} \mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] = \arg\max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi]$$

# Max Entropy RL

In maximum entropy RL, we change the objective

$$\arg \max_{\theta_\pi} \mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] = \arg \max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi]$$

We consider the entropy in the return

# Max Entropy RL

In maximum entropy RL, we change the objective

$$\arg\max_{\theta_\pi} \mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] = \arg\max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi\big]$$

We consider the entropy in the return

$$\arg\max_{\theta_\pi} \mathbb{E}[\mathcal{H}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] =$$

$$\arg\max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi\big] + \gamma^t H(\pi(a \mid s_t; \theta_\pi))$$

# Max Entropy RL

In maximum entropy RL, we change the objective

$$\arg\max_{\theta_\pi} \mathbb{E}[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] = \arg\max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi]$$

We consider the entropy in the return

$$\arg\max_{\theta_\pi} \mathbb{E}[\mathcal{H}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] =$$

$$\arg\max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi] + \gamma^t H(\pi(a \mid s_t; \theta_\pi))$$

We want a policy that is both random and maximizes the return

# Max Entropy RL

$$\arg\max_{\theta_\pi} \mathbb{E}[\mathcal{H}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] =$$

$$\arg\max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi] + H(\pi(\cdot \mid s_t; \theta_\pi)) \right)$$

# Max Entropy RL

$$\arg\max_{\theta_\pi} \mathbb{E}\big[\mathcal{H}(\boldsymbol{\tau}) \mid s_0; \theta_\pi\big] =$$

$$\arg\max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \big(\mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi\big] + H(\pi(\cdot \mid s_t; \theta_\pi))\big)$$

**Question:** Why do want policy entropy? Lowers the return

**Answer:** No good theoretical reason, helpful in practice

# Max Entropy RL

$$\arg \max_{\theta_\pi} \mathbb{E}\big[\mathcal{H}(\boldsymbol{\tau}) \mid s_0; \theta_\pi\big] =$$

$$\arg \max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \big(\mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi\big] + H(\pi(\cdot \mid s_t; \theta_\pi))\big)$$

**Question:** Why do want policy entropy? Lowers the return

**Answer:** No good theoretical reason, helpful in practice

- There are theoretical reasons, but they are not very good

# Max Entropy RL

$$\arg \max_{\theta_\pi} \mathbb{E}[\mathcal{H}(\tau) \mid s_0; \theta_\pi] =$$

$$\arg \max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi] + H(\pi(\cdot \mid s_t; \theta_\pi)) \right)$$

**Question:** Why do want policy entropy? Lowers the return

**Answer:** No good theoretical reason, helpful in practice
- There are theoretical reasons, but they are not very good
- Better exploration during training

# Max Entropy RL

$$\arg\max_{\theta_\pi} \mathbb{E}\big[\mathcal{H}(\boldsymbol{\tau}) \mid s_0; \theta_\pi\big] =$$

$$\arg\max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \big(\mathbb{E}\big[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi\big] + H(\pi(\cdot \mid s_t; \theta_\pi))\big)$$

**Question:** Why do want policy entropy? Lowers the return

**Answer:** No good theoretical reason, helpful in practice
- There are theoretical reasons, but they are not very good
- Better exploration during training
- More stable training (explain further in offline RL)

# Max Entropy RL

DDPG is based on a deterministic policy $\mu$

# Max Entropy RL

DDPG is based on a deterministic policy $\mu$

A deterministic policy $\mu$ has no entropy (it is not random)

# Max Entropy RL

DDPG is based on a deterministic policy $\mu$

A deterministic policy $\mu$ has no entropy (it is not random)

But for SAC, I talk about $\pi$ and policy entropy

# Max Entropy RL

DDPG is based on a deterministic policy $\mu$

A deterministic policy $\mu$ has no entropy (it is not random)

But for SAC, I talk about $\pi$ and policy entropy

How is this possible?

# Max Entropy RL

DDPG is based on a deterministic policy $\mu$

A deterministic policy $\mu$ has no entropy (it is not random)

But for SAC, I talk about $\pi$ and policy entropy

How is this possible?

Consider a function $f : S \times \Theta \times \mathbb{R} \mapsto A$

# Max Entropy RL

DDPG is based on a deterministic policy $\mu$

A deterministic policy $\mu$ has no entropy (it is not random)

But for SAC, I talk about $\pi$ and policy entropy

How is this possible?

Consider a function $f : S \times \Theta \times \mathbb{R} \mapsto A$

$$a = f\left(s, \theta_\mu, \eta\right) = \mu\left(s, \theta_\mu\right) + \eta$$

# Max Entropy RL

$$a = f(s, \theta_\mu, \eta) = \mu(s, \theta_\mu) + \eta$$

# Max Entropy RL

$$a = f(s, \theta_\mu, \eta) = \mu(s, \theta_\mu) + \eta$$

If $\eta \sim \text{Normal}(0, 1)$ our function $f$ is still deterministic

# Max Entropy RL

$$a = f(s, \theta_\mu, \eta) = \mu(s, \theta_\mu) + \eta$$

If $\eta \sim \text{Normal}(0, 1)$ our function $f$ is still deterministic

We call this the **reparameterization trick**, used in variational autoencoders (VAE)

# Max Entropy RL

$$a = f\left(s, \theta_\mu, \eta\right) = \mu\left(s, \theta_\mu\right) + \eta$$

If $\eta \sim \mathrm{Normal}(0, 1)$ our function $f$ is still deterministic

We call this the **reparameterization trick**, used in variational autoencoders (VAE)

Mathematically, this is a "deterministic" policy

# Max Entropy RL

$$a = f(s, \theta_\mu, \eta) = \mu(s, \theta_\mu) + \eta$$

If $\eta \sim \text{Normal}(0, 1)$ our function $f$ is still deterministic

We call this the **reparameterization trick**, used in variational autoencoders (VAE)

Mathematically, this is a "deterministic" policy

$$\mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0, \eta_0, \eta_1, ...; \theta_\mu\big]$$

# Max Entropy RL

$$a = f(s, \theta_\mu, \eta) = \mu(s, \theta_\mu) + \eta$$

$$\mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0, \eta_0, \eta_1, ...; \theta_\mu\big]$$

In practice, $f$ behaves just like a stochastic policy $\pi$

# Max Entropy RL

$$a = f\left(s, \theta_\mu, \eta\right) = \mu\left(s, \theta_\mu\right) + \eta$$

$$\mathbb{E}\left[\mathcal{G}(\boldsymbol{\tau}) \mid s_0, \eta_0, \eta_1, ...; \theta_\mu\right]$$

In practice, $f$ behaves just like a stochastic policy $\pi$

Gradient descent will learn $\theta_\mu$ that generalize over $\eta$

# Max Entropy RL

$$a = f\left(s, \theta_\mu, \eta\right) = \mu\left(s, \theta_\mu\right) + \eta$$

$$\mathbb{E}\left[\mathcal{G}(\boldsymbol{\tau}) \mid s_0, \eta_0, \eta_1, ...; \theta_\mu\right]$$

In practice, $f$ behaves just like a stochastic policy $\pi$

Gradient descent will learn $\theta_\mu$ that generalize over $\eta$

We abuse notation and write $f$ as a random policy

# Max Entropy RL

$$a = f\left(s, \theta_\mu, \eta\right) = \mu\left(s, \theta_\mu\right) + \eta$$

$$\mathbb{E}\left[\mathcal{G}(\boldsymbol{\tau}) \mid s_0, \eta_0, \eta_1, ...; \theta_\mu\right]$$

In practice, $f$ behaves just like a stochastic policy $\pi$

Gradient descent will learn $\theta_\mu$ that generalize over $\eta$

We abuse notation and write $f$ as a random policy

$$\pi(a \mid s; \theta_\pi)$$

# Max Entropy RL

$$a = f\left(s, \theta_\mu, \eta\right) = \mu\left(s, \theta_\mu\right) + \eta$$

$$\mathbb{E}\left[\mathcal{G}(\boldsymbol{\tau}) \mid s_0, \eta_0, \eta_1, ...; \theta_\mu\right]$$

In practice, $f$ behaves just like a stochastic policy $\pi$

Gradient descent will learn $\theta_\mu$ that generalize over $\eta$

We abuse notation and write $f$ as a random policy

$$\pi(a \mid s; \theta_\pi)$$

But $f$ is deterministic when we know $\eta$

# Max Entropy RL

What happens when we combine max entropy RL

# Max Entropy RL

What happens when we combine max entropy RL

$$\arg \max_{\theta_\pi} \mathbb{E}[\mathcal{H}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] =$$

$$\arg \max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi] + \gamma^t H(\pi(a \mid s_t; \theta_\pi))$$

# Max Entropy RL

What happens when we combine max entropy RL

$$\arg\max_{\theta_\pi} \mathbb{E}[\mathcal{H}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] =$$

$$\arg\max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi] + \gamma^t H(\pi(a \mid s_t; \theta_\pi))$$

With the reparameterization trick?

# Max Entropy RL

What happens when we combine max entropy RL

$$\arg \max_{\theta_\pi} \mathbb{E}[\mathcal{H}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] =$$

$$\arg \max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi] + \gamma^t H(\pi(a \mid s_t; \theta_\pi))$$

With the reparameterization trick?

$$a = \mu(s, \theta_\mu) + \eta \cong a \sim \pi(\cdot \mid s; \theta_\pi)$$

# Max Entropy RL

What happens when we combine max entropy RL

$$\arg \max_{\theta_\pi} \mathbb{E}[\mathcal{H}(\boldsymbol{\tau}) \mid s_0; \theta_\pi] =$$

$$\arg \max_{\theta_\pi} \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[\mathcal{R}(s_{t+1}) \mid s_0; \theta_\pi] + \gamma^t H(\pi(a \mid s_t; \theta_\pi))$$

With the reparameterization trick?

$$a = \mu(s, \theta_\mu) + \eta \cong a \sim \pi(\cdot \mid s; \theta_\pi)$$

We get SAC!

# Max Entropy RL

**Definition:** Soft Actor Critic (SAC) adds a max entropy objective and stochastic policy to DDPG

# Max Entropy RL

**Definition:** Soft Actor Critic (SAC) adds a max entropy objective and stochastic policy to DDPG

**Step 1:** Learn a $Q$ function for max entropy policy (Q learning)

# Max Entropy RL

**Definition:** Soft Actor Critic (SAC) adds a max entropy objective and stochastic policy to DDPG

**Step 1:** Learn a $Q$ function for max entropy policy (Q learning)

$$\theta_{Q,i+1} = \arg \min_{\theta_{Q,i}} \left( Q(s_0, a_0, \theta_{\pi,i}, \theta_{Q,i}) - y \right)^2$$

# Max Entropy RL

**Definition:** Soft Actor Critic (SAC) adds a max entropy objective and stochastic policy to DDPG

**Step 1:** Learn a $Q$ function for max entropy policy (Q learning)

$$\theta_{Q,i+1} = \arg\min_{\theta_{Q,i}} \left( Q\left(s_0, a_0, \theta_{\pi,i}, \theta_{Q,i}\right) - y \right)^2$$

$$y = \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + H\left(\pi\left(a \mid s_0; \theta_\mu\right)\right) + \gamma Q\left(s_1, \mu(s_1, \theta_\mu, \eta), \theta_{\pi,i}, \theta_{Q,i}\right)$$

# Max Entropy RL

**Definition:** Soft Actor Critic (SAC) adds a max entropy objective and stochastic policy to DDPG

**Step 1:** Learn a $Q$ function for max entropy policy (Q learning)

$$\theta_{Q,i+1} = \arg \min_{\theta_{Q,i}} \left( Q\left(s_0, a_0, \theta_{\pi,i}, \theta_{Q,i}\right) - y\right)^2$$

$$y = \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + H\big(\pi(a \mid s_0; \theta_\mu)\big) + \gamma Q\big(s_1, \mu(s_1, \theta_\mu, \eta), \theta_{\pi,i}, \theta_{Q,i}\big)$$

# Max Entropy RL

**Definition:** Soft Actor Critic (SAC) adds a max entropy objective and stochastic policy to DDPG

**Step 1:** Learn a $Q$ function for max entropy policy (Q learning)

$$\theta_{Q,i+1} = \arg\min_{\theta_{Q,i}} \left( Q\big(s_0, a_0, \theta_{\pi,i}, \theta_{Q,i}\big) - y \right)^2$$

$$y = \underbrace{\hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0]}_{\text{Reward}} + H\big(\pi\big(a \mid s_0; \theta_\mu\big)\big) + \gamma Q\big(s_1, \mu(s_1, \theta_\mu, \eta), \theta_{\pi,i}, \theta_{Q,i}\big)$$

# Max Entropy RL

**Definition:** Soft Actor Critic (SAC) adds a max entropy objective and stochastic policy to DDPG

**Step 1:** Learn a $Q$ function for max entropy policy (Q learning)

$$\theta_{Q,i+1} = \arg \min_{\theta_{Q,i}} \left( Q\big(s_0, a_0, \theta_{\pi,i}, \theta_{Q,i}\big) - y \right)^2$$

$$y = \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + H\big(\pi\big(a \mid s_0; \theta_\mu\big)\big) + \gamma Q\big(s_1, \mu(s_1, \theta_\mu, \eta), \theta_{\pi,i}, \theta_{Q,i}\big)$$

Reward

Entropy bonus

# Max Entropy RL

**Definition:** Soft Actor Critic (SAC) adds a max entropy objective and stochastic policy to DDPG

**Step 1:** Learn a $Q$ function for max entropy policy (Q learning)

$$\theta_{Q,i+1} = \arg \min_{\theta_{Q,i}} \left( Q\big(s_0, a_0, \theta_{\pi,i}, \theta_{Q,i}\big) - y\right)^2$$

$$y = \hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0] + H\big(\pi\big(a \mid s_0; \theta_\mu\big)\big) + \gamma Q\big(s_1, \mu(s_1, \theta_\mu, \eta), \theta_{\pi,i}, \theta_{Q,i}\big)$$

Reward      Entropy bonus      Deterministic $a$

# Max Entropy RL

**Definition:** Soft Actor Critic (SAC) adds a max entropy objective and stochastic policy to DDPG

**Step 1:** Learn a $Q$ function for max entropy policy (Q learning)

$$\theta_{Q,i+1} = \arg \min_{\theta_{Q,i}} \left( Q(s_0, a_0, \theta_{\pi,i}, \theta_{Q,i}) - y \right)^2$$

$$y = \underbrace{\hat{\mathbb{E}}[\mathcal{R}(s_1) \mid s_0, a_0]}_{\text{Reward}} + \underbrace{H(\pi(a \mid s_0; \theta_\mu))}_{\text{Entropy bonus}} + \gamma Q\big(s_1, \underbrace{\mu(s_1, \theta_\mu, \eta)}_{\text{Deterministic } a}, \theta_{\pi,i}, \theta_{Q,i}\big)$$

where $\eta$ is randomly sampled

# Max Entropy RL

**Step 2:** Learn a $\pi$ that maximizes $Q$ (policy gradient)

# Max Entropy RL

**Step 2:** Learn a $\pi$ that maximizes $Q$ (policy gradient)

$$\theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \underbrace{Q\big(s_0, \mu(s_0, \theta_\mu, \eta), \theta_\pi, \theta_Q\big)}_{\text{Replaces } \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]}$$

# Max Entropy RL

**Step 2:** Learn a $\pi$ that maximizes $Q$ (policy gradient)

Deterministic $a$

$$\theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \underbrace{Q\big(s_0, \mu(s_0, \theta_\mu, \eta), \theta_\pi, \theta_Q\big)}_{\text{Replaces } \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]}$$

# Max Entropy RL

**Step 2:** Learn a $\pi$ that maximizes $Q$ (policy gradient)

$$\theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \underbrace{Q\big(s_0, \overset{\text{Deterministic } a}{\mu(s_0, \theta_\mu, \eta)}, \theta_\pi, \theta_Q\big)}_{\text{Replaces } \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]}$$

where $\eta$ is randomly sampled

# Max Entropy RL

**Step 2:** Learn a $\pi$ that maximizes $Q$ (policy gradient)

$$\theta_{\pi,i+1} = \theta_{\pi,i} + \alpha \cdot \underbrace{Q\big(s_0, \overbrace{\mu(s_0, \theta_\mu, \eta)}^{\text{Deterministic } a}, \theta_\pi, \theta_Q\big)}_{\text{Replaces } \mathbb{E}\big[\mathcal{G}(\boldsymbol{\tau}) \mid s_0; \theta_\mu\big]}$$

where $\eta$ is randomly sampled

Repeat until convergence, $\theta_{\mu,i+1} = \theta_{\mu,i}, \quad \theta_{Q,i+1} = \theta_{Q,i}$

# Max Entropy RL

Like PPO, there are many variants of SAC

# Max Entropy RL

Like PPO, there are many variants of SAC

- Learn separate value and Q functions

# Max Entropy RL

Like PPO, there are many variants of SAC

- Learn separate value and Q functions
- Double Q function

# Max Entropy RL

Like PPO, there are many variants of SAC

- Learn separate value and Q functions
- Double Q function
- Lagrangian entropy adjustment

# Max Entropy RL

Like PPO, there are many variants of SAC

- Learn separate value and Q functions
- Double Q function
- Lagrangian entropy adjustment
- Reduced variance gradients

# Max Entropy RL

Like PPO, there are many variants of SAC

- Learn separate value and Q functions
- Double Q function
- Lagrangian entropy adjustment
- Reduced variance gradients

Like PPO, SAC is complicated – uses many "implementation tricks"

# Max Entropy RL

Like PPO, there are many variants of SAC

- Learn separate value and Q functions
- Double Q function
- Lagrangian entropy adjustment
- Reduced variance gradients

Like PPO, SAC is complicated – uses many "implementation tricks"

- Often not documented
- CleanRL describes modern SAC, using tricks from 5+ papers
- https://docs.cleanrl.dev/rl-algorithms/sac/#implementation-details_1

Coding SAC could take an entire lecture, read CleanRL

# Max Entropy RL

DDPG $\approx$ A2C

# Max Entropy RL

DDPG $\approx$ A2C

Introduces the concept, simple

# Max Entropy RL

DDPG ≈ A2C

Introduces the concept, simple

SAC ≈ PPO

# Max Entropy RL

DDPG $\approx$ A2C

Introduces the concept, simple

SAC $\approx$ PPO

Improves the concept, complex

# Max Entropy RL

DDPG $\approx$ A2C

SAC $\approx$ PPO

Introduces the concept, simple

Improves the concept, complex

I suggest you try DDPG before SAC

Actor Critic II

# Max Entropy RL

DDPG $\approx$ A2C

SAC $\approx$ PPO

Introduces the concept, simple

Improves the concept, complex

I suggest you try DDPG before SAC
- DDPG is much easier to implement

# Max Entropy RL

DDPG ≈ A2C                          SAC ≈ PPO

Introduces the concept, simple     Improves the concept, complex

I suggest you try DDPG before SAC
- DDPG is much easier to implement
- Fewer hyperparameters

# Max Entropy RL

DDPG ≈ A2C

SAC ≈ PPO

Introduces the concept, simple

Improves the concept, complex

I suggest you try DDPG before SAC

- DDPG is much easier to implement
- Fewer hyperparameters
- Tuned DDPG can likely outperform untuned SAC

# Max Entropy RL

What algorithm is best in 2025?

# Max Entropy RL

What algorithm is best in 2025?

For discrete actions (Atari), DQN variants still perform best

# Max Entropy RL

What algorithm is best in 2025?

For discrete actions (Atari), DQN variants still perform best

For continuous actions (MuJoCo), SAC performs best

# Max Entropy RL

What algorithm is best in 2025?

For discrete actions (Atari), DQN variants still perform best

For continuous actions (MuJoCo), SAC performs best

PPO is less sensitive to hyperparameters than SAC/DQN

# Max Entropy RL

What algorithm is best in 2025?

For discrete actions (Atari), DQN variants still perform best

For continuous actions (MuJoCo), SAC performs best

PPO is less sensitive to hyperparameters than SAC/DQN
- Noobs like PPO, they cannot tune hyperparameters

# Max Entropy RL

What algorithm is best in 2025?

For discrete actions (Atari), DQN variants still perform best

For continuous actions (MuJoCo), SAC performs best

PPO is less sensitive to hyperparameters than SAC/DQN
- Noobs like PPO, they cannot tune hyperparameters
- Often performs worse than tuned DQN/SAC

# Max Entropy RL

What algorithm is best in 2025?

For discrete actions (Atari), DQN variants still perform best

For continuous actions (MuJoCo), SAC performs best

PPO is less sensitive to hyperparameters than SAC/DQN
- Noobs like PPO, they cannot tune hyperparameters
- Often performs worse than tuned DQN/SAC

Tuned DDPG/A2C perform 95% as good as tuned SAC/PPO

# Max Entropy RL

What algorithm is best in 2025?

For discrete actions (Atari), DQN variants still perform best

For continuous actions (MuJoCo), SAC performs best

PPO is less sensitive to hyperparameters than SAC/DQN
- Noobs like PPO, they cannot tune hyperparameters
- Often performs worse than tuned DQN/SAC

Tuned DDPG/A2C perform 95% as good as tuned SAC/PPO
- Much easier to debug

# Final Project Tips

# Final Project Tips

Log and plot EVERYTHING

Actor Critic II

# Final Project Tips

Log and plot EVERYTHING

- Losses, mean advantage, mean Q, policy entropy, etc
    - ▸ Use these to help debug and tune hyperparameters
    - ▸ E.g., exploding losses, decrease learning rate
    - ▸ E.g., Q values too large? Increase time between target net update

# Final Project Tips

Log and plot EVERYTHING

- Losses, mean advantage, mean Q, policy entropy, etc
  - ‣ Use these to help debug and tune hyperparameters
  - ‣ E.g., exploding losses, decrease learning rate
  - ‣ E.g., Q values too large? Increase time between target net update

If you get stuck, visualize your policy

# Final Project Tips

Log and plot EVERYTHING

- Losses, mean advantage, mean Q, policy entropy, etc
  - ‣ Use these to help debug and tune hyperparameters
  - ‣ E.g., exploding losses, decrease learning rate
  - ‣ E.g., Q values too large? Increase time between target net update

If you get stuck, visualize your policy

- Record some episodes (videos/frames/etc)

# Final Project Tips

Log and plot EVERYTHING

- Losses, mean advantage, mean Q, policy entropy, etc
  - ▸ Use these to help debug and tune hyperparameters
  - ▸ E.g., exploding losses, decrease learning rate
  - ▸ E.g., Q values too large? Increase time between target net update

If you get stuck, visualize your policy

- Record some episodes (videos/frames/etc)
- Watch the policy, what did it learn to do?

# Final Project Tips

Log and plot EVERYTHING

- Losses, mean advantage, mean Q, policy entropy, etc
  - ‣ Use these to help debug and tune hyperparameters
  - ‣ E.g., exploding losses, decrease learning rate
  - ‣ E.g., Q values too large? Increase time between target net update

If you get stuck, visualize your policy

- Record some episodes (videos/frames/etc)
- Watch the policy, what did it learn to do?
  - ‣ Think about why it learned to do this (exploiting bugs in MDP)

# Final Project Tips

Why does Steven spend so much time on theory instead of coding?

# Final Project Tips

Why does Steven spend so much time on theory instead of coding?

In supervised learning, follow MNIST tutorial and everything works

# Final Project Tips

Why does Steven spend so much time on theory instead of coding?

In supervised learning, follow MNIST tutorial and everything works

This is **NOT** the case in RL

# Final Project Tips

Why does Steven spend so much time on theory instead of coding?

In supervised learning, follow MNIST tutorial and everything works

This is **NOT** the case in RL

In RL, your code never works on the first try

# Final Project Tips

Why does Steven spend so much time on theory instead of coding?

In supervised learning, follow MNIST tutorial and everything works

This is **NOT** the case in RL

In RL, your code never works on the first try

Even if it is correct, you need to play with hyperparameters

# Final Project Tips

Why does Steven spend so much time on theory instead of coding?

In supervised learning, follow MNIST tutorial and everything works

This is **NOT** the case in RL

In RL, your code never works on the first try

Even if it is correct, you need to play with hyperparameters

Theory is absolutely necessary to understand **why** your policy fails, and **how** to fix it

# Final Project Tips

Why does Steven spend so much time on theory instead of coding?

In supervised learning, follow MNIST tutorial and everything works

This is **NOT** the case in RL

In RL, your code never works on the first try

Even if it is correct, you need to play with hyperparameters

Theory is absolutely necessary to understand **why** your policy fails, and **how** to fix it

You must use your brain to be successful!