

CS4400

DEEP REINFORCEMENT LEARNING

Lecture 7: Off-Policy Actor-Critic

Wendelin Böhmer

`<j.w.bohmer@tudelft.nl>`



12th of December 2023

Content of this lecture



- 7.1 Deterministic policies
- 7.2 Robust decision making
- 7.3 Maximum entropy RL

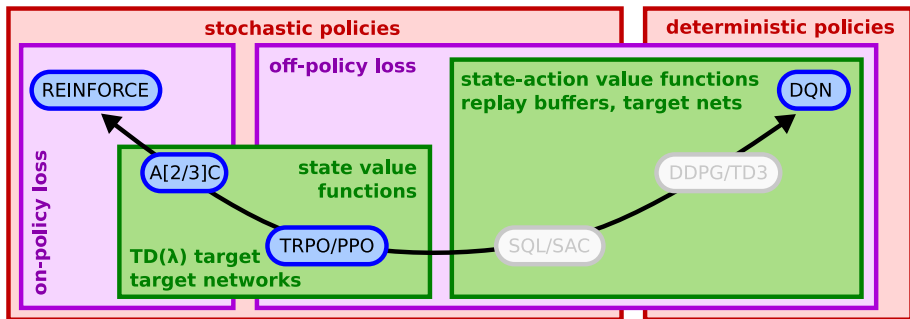
7.1

Off-Policy Actor-Critic Deterministic policies

7.1 Deep RL algorithms so far



- Spectrum from REINFORCE to DQN
 - Algorithm choice depends how well value functions generalize
- Empirically PPO best actor-critic algorithms so far
 - allows off-policy training on batches of on-policy samples





core concept: Deterministic policy gradients



- Continuous actions $\mathbf{a} \in \mathcal{A} \subset \mathbb{R}^m$

$$\nabla_{\theta} \mathcal{L}_{\pi}[\theta] \approx -\nabla_{\theta} \mathbb{E}_{\mu} \left[\sum_{t=0}^{n-1} \gamma^t \left(\overbrace{Q^{\pi}(s_t, \mathbf{a}_t)}^{A^{\pi}(s_t, \mathbf{a}_t)} - V^{\pi}(s_t) \right) \frac{\pi_{\theta}(\mathbf{a}_t | s_t)}{\mu(\mathbf{a}_t | s_t)} \right]$$

approximation drops $\frac{\xi_t^{\pi}(s_t)}{\xi^{\mu}(s_t)}$, see [Sutton et al. \(2016\)](#) for ways to estimate it

(DPG, [Silver et al., 2014](#))



core concept: Deterministic policy gradients



- Continuous actions $\mathbf{a} \in \mathcal{A} \subset \mathbb{R}^m$

$$\begin{aligned}\nabla_{\theta} \mathcal{L}_{\pi}[\theta] &\approx -\nabla_{\theta} \mathbb{E}_{\mu} \left[\sum_{t=0}^{n-1} \gamma^t \left(\overbrace{Q^{\pi}(s_t, \mathbf{a}_t)}^{A^{\pi}(s_t, \mathbf{a}_t)} - V^{\pi}(s_t) \right) \frac{\pi_{\theta}(\mathbf{a}_t | s_t)}{\mu(\mathbf{a}_t | s_t)} \right] \\ &= -\nabla_{\theta} \sum_{t=0}^{n-1} \iint \xi_t^{\mu}(s_t) \pi_{\theta}(\mathbf{a}_t | s_t) \gamma^t Q^{\pi}(s_t, \mathbf{a}_t) ds_t d\mathbf{a}_t\end{aligned}$$

approximation drops $\frac{\xi_t^{\pi}(s_t)}{\xi^{\mu}(s_t)}$, see [Sutton et al. \(2016\)](#) for ways to estimate it

(DPG, [Silver et al., 2014](#))



core concept: Deterministic policy gradients



- Continuous actions $\mathbf{a} \in \mathcal{A} \subset \mathbb{R}^m$
- Deterministic policy $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$
 - equivalent to $\pi_\theta(\pi_\theta(s)|s) = 1, \forall s \in \mathcal{S}$

$$\begin{aligned}\nabla_\theta \mathcal{L}_\pi[\theta] &\approx -\nabla_\theta \mathbb{E}_\mu \left[\sum_{t=0}^{n-1} \gamma^t \left(\overbrace{Q^\pi(s_t, \mathbf{a}_t)}^{A^\pi(s_t, \mathbf{a}_t)} - V^\pi(s_t) \right) \frac{\pi_\theta(\mathbf{a}_t|s_t)}{\mu(\mathbf{a}_t|s_t)} \right] \\ &= -\nabla_\theta \sum_{t=0}^{n-1} \iint \xi_t^\mu(s_t) \pi_\theta(\mathbf{a}_t|s_t) \gamma^t Q^\pi(s_t, \mathbf{a}_t) ds_t d\mathbf{a}_t \\ &= -\sum_{t=0}^{n-1} \int \xi_t^\mu(s_t) \gamma^t \nabla_\theta Q^\pi(s_t, \pi_\theta(s_t)) ds_t\end{aligned}$$

approximation drops $\frac{\xi_t^\pi(s_t)}{\xi_t^\mu(s_t)}$, see [Sutton et al. \(2016\)](#) for ways to estimate it

(DPG, [Silver et al., 2014](#))



core concept: Deterministic policy gradients



- Continuous actions $\mathbf{a} \in \mathcal{A} \subset \mathbb{R}^m$
- Deterministic policy $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$
 - equivalent to $\pi_\theta(\pi_\theta(s)|s) = 1, \forall s \in \mathcal{S}$

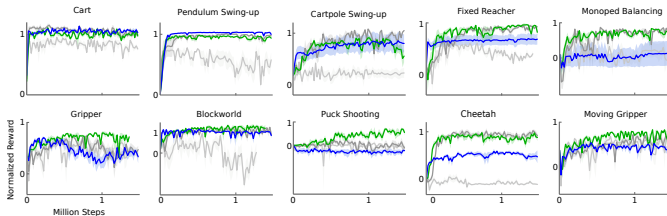
$$\begin{aligned}\nabla_\theta \mathcal{L}_\pi[\theta] &\approx -\nabla_\theta \mathbb{E}_\mu \left[\sum_{t=0}^{n-1} \gamma^t \left(\overbrace{Q^\pi(s_t, \mathbf{a}_t)}^{A^\pi(s_t, \mathbf{a}_t)} - V^\pi(s_t) \right) \frac{\pi_\theta(\mathbf{a}_t|s_t)}{\mu(\mathbf{a}_t|s_t)} \right] \\&= -\nabla_\theta \sum_{t=0}^{n-1} \iint \xi_t^\mu(s_t) \pi_\theta(\mathbf{a}_t|s_t) \gamma^t Q^\pi(s_t, \mathbf{a}_t) ds_t d\mathbf{a}_t \\&= -\sum_{t=0}^{n-1} \int \xi_t^\mu(s_t) \gamma^t \nabla_\theta Q^\pi(s_t, \pi_\theta(s_t)) ds_t \\&\propto -\nabla_\theta \mathbb{E}_\mu \left[Q^\pi(s_t, \pi_\theta(s_t)) \middle| t \sim p(\cdot|\gamma) \right], \quad p(t|\gamma) \propto \gamma^t\end{aligned}$$

- Implementations based on transitions often omit γ^t

$p(t|\gamma)$ is technically a [Conway-Maxwell-Poisson distribution](#) with $\lambda = \gamma$ and $\nu = 0$.

(DPG, [Silver et al., 2014](#))

- Off-policy DPG with:
 - Experience replay buffer
 - Actor and critic target networks (TN)
 - Soft target update
 - Batch-norm over inputs (BN)
 - Exploration by adding noise to $\pi_{\theta}(s_t)$



DPG + BN
 DPG + TN
DPG + TN + BN
DPG + TN (pixels)

spinningup.openai.com/en/latest/algorithms/ddpg.html

DDPG and images/results by [Lillicrap et al. \(2016\)](#)



Question: Implementing DPG



- For actor network $\pi_{\theta}(s)$ and critic network $Q_{\phi}(s, \mathbf{a})$:

$$\nabla_{\theta} \mathcal{L}_{\pi}[\theta] \approx - \sum_{t=0}^{n-1} \gamma^t \mathbb{E}_{\mu} \left[\nabla_{\theta} Q_{\phi}(s_t, \pi_{\theta}(s_t)) \right]$$

- How can we compute the gradients for θ and ϕ in PyTorch?



Question: Implementing DPG

- For actor network $\pi_{\theta}(s)$ and critic network $Q_{\phi}(s, \mathbf{a})$:

$$\nabla_{\theta} \mathcal{L}_{\pi}[\theta] \approx - \sum_{t=0}^{n-1} \gamma^t \mathbb{E}_{\mu} \left[\nabla_{\theta} Q_{\phi}(s_t, \pi_{\theta}(s_t)) \right]$$

- How can we compute the gradients for θ and ϕ in PyTorch?
- each network needs its own optimizer

$$\arg \min_{\phi} \mathbb{E}_{\mu} \left[\sum_{t=0}^{n-1} \left(r_t + \gamma Q_{\phi'}(s_{t+1}, \pi_{\theta'}(s_{t+1})) - Q_{\phi}(s_t, a_t) \right)^2 \right]$$

$$\arg \min_{\theta} - \mathbb{E}_{\mu} \left[\sum_{t=0}^{n-1} \gamma^t Q_{\phi}(s_t, \pi_{\theta}(s_t)) \right]$$

- Try to write the above gradient descend problems in PyTorch!
 - given a mini-batch `batch = self.replay_buffer.sample()`
 - given an actor (`batch['states']`)
 - given a critic (`batch['states']`, `batch['actions']`)

7.1 DDPG implementation

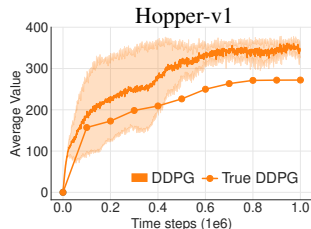


```
1 from torch.optim import Adam
2 from torch.nn.functional import mse_loss
3 actor_optimizer = Adam(actor.parameters())
4 critic_optimizer = Adam(critic.parameters())
5 actor_t, critic_t = deepcopy(actor), deepcopy(critic)
6 for _ in range(max_updates):
7     batch = self.replay_buffer.sample()
8     # Critic update
9     q_values = critic(batch['states'], batch['actions'])
10    next_actions = actor_t(batch['next_states'])
11    targets = batch['rewards'] + gamma * (~batch['terminals'] \
12        * critic_t(batch['next_states'], next_actions))
13    critic_optimizer.zero_grad()
14    mse_loss(q_values, targets.detach()).backward()
15    critic_optimizer.step()
16    # Actor update
17    q_values = critic(batch['states'], actor(batch['states']))
18    actor_optimizer.zero_grad()
19    (-q_values * batch['discounts']).mean().backward()
20    actor_optimizer.step()
21    # Target network updates
22    actor_t = target_model_updates(actor_t, actor)
23    critic_t = target_model_updates(critic_t, critic)
```

7.1 Twin delayed DDPG (TD3)



- $\mathbb{E}[\max_{\theta} Q(s, \pi_{\theta}(s))] \geq \max_{\theta} \mathbb{E}[Q(s, \pi_{\theta}(s))]$
 - like in Double Q-learning (Lecture 5)

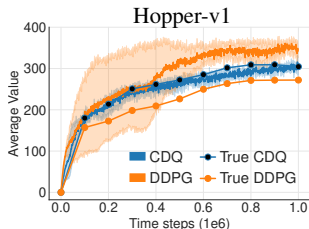


images and results from [Fujimoto et al. \(TD3, 2018\)](#)

7.1 Twin delayed DDPG (TD3)



- $\mathbb{E}[\max_{\theta} Q(s, \pi_{\theta}(s))] \geq \max_{\theta} \mathbb{E}[Q(s, \pi_{\theta}(s))]$
 - like in Double Q-learning (Lecture 5)
- CDQ: two critics against overestimation
 - minimum, *unlike* in Double Q-learning!



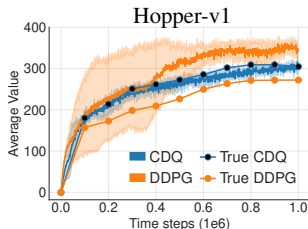
$$\mathcal{L}_i^Q[\phi_i] = \mathbb{E}_{\mu} \left[\sum_{t=0}^{n-1} \left(r_t + \gamma \min_{j \in \{1,2\}} Q_{\phi'_j}(s_{t+1}, \pi_{\theta'}(s_{t+1})) - Q_{\phi_i}(s_t, a_t) \right)^2 \right]$$

Clipped Double Q-learning (CDQ) is Q-learning with two Q-value functions images and results from [Fujimoto et al. \(TD3, 2018\)](#)

7.1 Twin delayed DDPG (TD3)



- $\mathbb{E}[\max_{\theta} Q(s, \pi_{\theta}(s))] \geq \max_{\theta} \mathbb{E}[Q(s, \pi_{\theta}(s))]$
 - like in Double Q-learning (Lecture 5)
- CDQ: two critics against overestimation
 - minimum, *unlike* in Double Q-learning!
- TD3: regularization with clipped noise
 $\epsilon \sim \text{clip}(\mathcal{N}(\mathbf{0}, \sigma^2), -c, c) \in \mathbb{R}^{|A|}$



$$\mathcal{L}_i^Q[\phi_i] = \mathbb{E}_{\mu} \left[\sum_{t=0}^{n-1} \left(r_t + \gamma \min_{j \in \{1,2\}} Q_{\phi'_j}(s_{t+1}, \pi_{\theta'}(s_{t+1}) + \epsilon) - Q_{\phi_i}(s_t, a_t) \right)^2 \right]$$

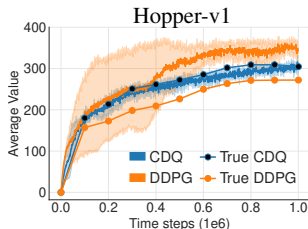
- TD3 algorithm one of the highest performing DDPG variants
 - introduces pessimistic values (\rightarrow Lecture 9)

7.1 Twin delayed DDPG (TD3)

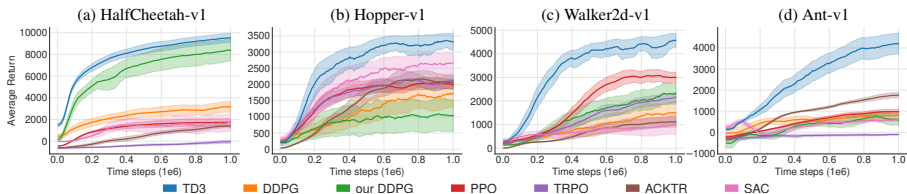


- $\mathbb{E}[\max_{\theta} Q(s, \pi_{\theta}(s))] \geq \max_{\theta} \mathbb{E}[Q(s, \pi_{\theta}(s))]$
 - like in Double Q-learning (Lecture 5)
- CDQ: two critics against overestimation
 - minimum, *unlike* in Double Q-learning!
- TD3: regularization with clipped noise

$$\epsilon \sim \text{clip}(\mathcal{N}(\mathbf{0}, \sigma^2), -c, c) \in \mathbb{R}^{|\mathcal{A}|}$$



$$\mathcal{L}_i^Q[\phi_i] = \mathbb{E}_{\mu} \left[\sum_{t=0}^{n-1} \left(r_t + \gamma \min_{j \in \{1,2\}} Q_{\phi'_j}(s_{t+1}, \pi_{\theta'}(s_{t+1}) + \epsilon) - Q_{\phi_i}(s_t, a_t) \right)^2 \right]$$

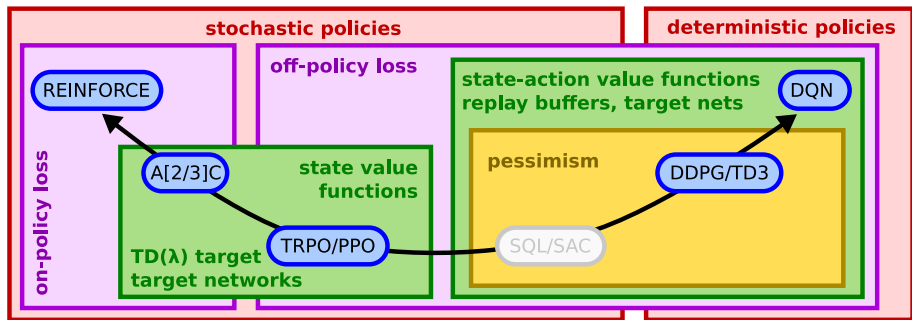


images and results from [Fujimoto et al. \(TD3, 2018\)](#)

7.1 DDPG/TD3 in comparison



- Only for deterministic policies and continuous actions
 - allows sample efficient off-policy sampling from replay buffer
 - state-action Q-value functions generalize worse
 - TD3 stabilizes Q-value with pessimism



TD3 is one of the highest performing algorithms for continuous actions



- DPG learns deterministic policies for continuous actions
- Requires Q-value of current policy function as critic
- Implementation with two optimizers
- TD3 is a modern version based on pessimism and action noise

Learning Objectives

LO7.1: Derive and implement the DPG losses

LO7.2: Explain the difference between DDPG and TD3

7.2

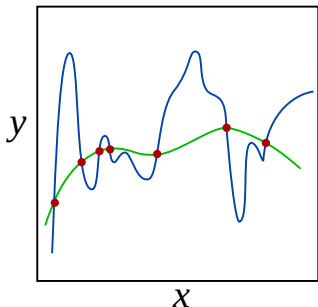
Off-Policy Actor-Critic

Robust decision making

7.2 Regularization



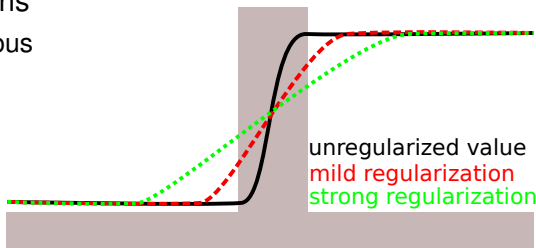
- TD3 used noisy actions as regularization of the policy
- Regularization smoothes functions by lowering ℓ , e.g. $\|a\|$
 - or averages over ensemble like Dropout
 - data augmentation sometimes called regularization



For regularization and generalization see Lecture 4

image source: www.wikipedia.org

- TD3 used noisy actions as regularization of the policy
- Regularization smoothes functions by lowering ℓ , e.g. $\|a\|$
 - or averages over ensemble like Dropout
 - data augmentation sometimes called regularization
- No effect for value functions
 - smooth values dangerous
 - e.g. near walls

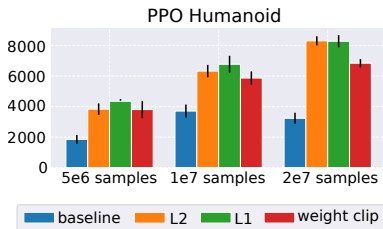


effect of regularization still open question, see [Liu et al. \(2021\)](#), image source) for details on presented example

7.2 Regularization



- TD3 used noisy actions as regularization of the policy
- Regularization smoothes functions by lowering ℓ , e.g. $\|a\|$
 - or averages over ensemble like Dropout
 - data augmentation sometimes called regularization
- No effect for value functions
 - smooth values dangerous
 - e.g. near walls
- Measurable effect for policies
 - why does this help?

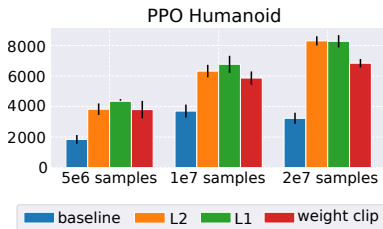


effect of regularization still open question, see [Liu et al. \(2021\)](#), image source) for details on presented example

7.2 Regularization



- TD3 used noisy actions as regularization of the policy
- Regularization smoothes functions by lowering ℓ , e.g. $\|a\|$
 - or averages over ensemble like Dropout
 - data augmentation sometimes called regularization
- No effect for value functions
 - smooth values dangerous
 - e.g. near walls
- Measurable effect for policies
 - why does this help?
 - regularized = smooth = stochastic
 - why are stochastic policies good?



effect of regularization still open question, see [Liu et al. \(2021\)](#), image source) for details on presented example

7.2 Robust reinforcement learning



- Approximation introduces local errors in the policy
 - online sampling corrects mistakes slowly
 - long-term consequences of small errors varies
 - robust decisions should have robust consequences
 - can improve learning even in deterministic MDP

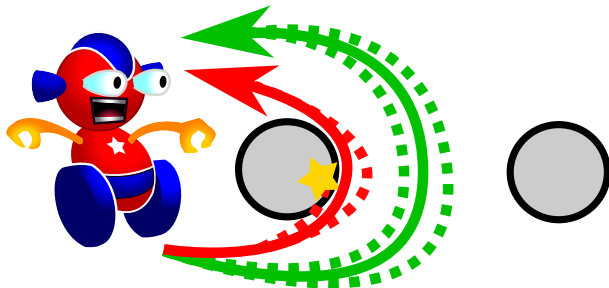


image source freesvg.org

- Approximation introduces local errors in the policy
 - online sampling corrects mistakes slowly
 - long-term consequences of small errors varies
 - robust decisions should have robust consequences
 - can improve learning even in deterministic MDP
- Robust RL disrupts MDP with term $\delta \in \Delta$
 - by changing transitions or rewards, e.g. $P(s'|s, a, \delta) := P(s' + \delta|s, a)$
 - withstanding disruption δ is rewarded, e.g. $\bar{r}(\delta) := \alpha \|\delta\|$

$$V_R^*(s) := \max_{a \in \mathcal{A}} \min_{\delta \in \Delta} \left(r(s, a, \delta) + \bar{r}(\delta) + \gamma \int P(s'|s, a, \delta) V_R^*(s') ds' \right)$$

basic concept by [Morimoto and Doya \(2001\)](#), see [slides](#) by Marek Petrik and the recent review by [Moos et al. \(2022\)](#) for details

- Approximation introduces local errors in the policy
 - online sampling corrects mistakes slowly
 - long-term consequences of small errors varies
 - robust decisions should have robust consequences
 - can improve learning even in deterministic MDP
 - Robust RL disrupts MDP with term $\delta \in \Delta$
 - by changing transitions or rewards, e.g. $P(s'|s, a, \delta) := P(s' + \delta | s, a)$
 - withstanding disruption δ is rewarded, e.g. $\bar{r}(\delta) := \alpha \|\delta\|$
- $$V_R^*(s) := \max_{a \in \mathcal{A}} \min_{\delta \in \Delta} \left(r(s, a, \delta) + \bar{r}(\delta) + \gamma \int P(s'|s, a, \delta) V_R^*(s') ds' \right)$$
- Requires controllable amounts of noise/disruptions δ
 - can we motivate $\bar{r}(\delta)$ without such control?
 - e.g. how much noise should be added to actions in TD3?

basic concept by [Morimoto and Doya \(2001\)](#), see [slides](#) by Marek Petrik and the recent review by [Moos et al. \(2022\)](#) for details

7.2 Planning as inference I



- Probability of trajectory $\tau_n = [s_0, a_0, \dots, s_n]$ under prior policy μ

$$p^\mu(\tau_n) = \rho(s_0) \prod_{t=0}^{n-1} \mu(a_t|s_t) P(s_{t+1}|s_t, a_t)$$

(PAI, Levine, 2018)

- Probability of trajectory $\tau_n = [s_0, a_0, \dots, s_n]$ under prior policy μ
 - introducing “optimality variables” $O_t \in \{0, 1\}$

$$p^\mu(\tau_n) = \rho(s_0) \prod_{t=0}^{n-1} \mu(a_t|s_t) P(s_{t+1}|s_t, a_t)$$

$$p^\mu(\tau_n, \{O_t\}_{t=0}^{n-1}) = \rho(s_0) \prod_{t=0}^{n-1} \mu(a_t|s_t) p(O_t|s_t, a_t) P(s_{t+1}|s_t, a_t)$$

- Probability of trajectory $\tau_n = [s_0, a_0, \dots, s_n]$ under prior policy μ
 - introducing “optimality variables” $O_t \in \{0, 1\}$
 - probability that all O_t are “optimal”

$$p^\mu(\tau_n) = \rho(s_0) \prod_{t=0}^{n-1} \mu(a_t|s_t) P(s_{t+1}|s_t, a_t)$$

$$p^\mu(\tau_n, \{O_t\}_{t=0}^{n-1}) = \rho(s_0) \prod_{t=0}^{n-1} \mu(a_t|s_t) p(O_t|s_t, a_t) P(s_{t+1}|s_t, a_t)$$

$$p^\mu(\tau_n | \{O_t=1\}_{t=0}^{n-1}) \propto \rho(s_0) \prod_{t=0}^{n-1} \mu(a_t|s_t) p(O_t=1|s_t, a_t) P(s_{t+1}|s_t, a_t)$$

- Probability of trajectory $\tau_n = [s_0, a_0, \dots, s_n]$ under prior policy μ
 - introducing “optimality variables” $O_t \in \{0, 1\}$
 - probability that all O_t are “optimal”
 - define optimality $p(O_t=1|s_t, a_t) \propto \exp(\frac{1}{\alpha}r(s_t, a_t))$

$$p^\mu(\tau_n) = \rho(s_0) \prod_{t=0}^{n-1} \mu(a_t|s_t) P(s_{t+1}|s_t, a_t)$$

$$p^\mu(\tau_n, \{O_t\}_{t=0}^{n-1}) = \rho(s_0) \prod_{t=0}^{n-1} \mu(a_t|s_t) p(O_t|s_t, a_t) P(s_{t+1}|s_t, a_t)$$

$$p^\mu(\tau_n | \{O_t=1\}_{t=0}^{n-1}) \propto \rho(s_0) \prod_{t=0}^{n-1} \mu(a_t|s_t) p(O_t=1|s_t, a_t) P(s_{t+1}|s_t, a_t)$$

$$p^\mu(\tau_n | \{O_t=1\}_{t=0}^{n-1}) \propto \rho(s_0) \prod_{t=0}^{n-1} \mu(a_t|s_t) \exp(\frac{1}{\alpha}r(s_t, a_t)) P(s_{t+1}|s_t, a_t)$$

7.2 Planning as inference II



- $p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})$ is the distribution of “optimal” trajectories!
 - Minimize KL-divergence between $p^{\pi_\theta}(\tau_n)$ and $p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})$

$$\mathcal{L}_\pi^{\text{PAI}} := D_{\text{KL}}[p^{\pi_\theta}(\tau_n) \parallel p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})]$$

(PAI, Levine, 2018)

- $p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})$ is the distribution of “optimal” trajectories!
 - Minimize KL-divergence between $p^{\pi_\theta}(\tau_n)$ and $p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})$

$$\begin{aligned} \mathcal{L}_\pi^{\text{PAI}} &:= D_{\text{KL}}[p^{\pi_\theta}(\tau_n) \parallel p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})] \\ &\propto \int p^{\pi_\theta}(\tau_n) \ln \frac{\rho(s_0) \prod_{t=0}^{n-1} \pi_\theta(a_t | s_t) P(s_{t+1} | s_t, a_t)}{\rho(s_0) \prod_{t=0}^{n-1} \mu(a_t | s_t) P(s_{t+1} | s_t, a_t) \exp(\frac{1}{\alpha} r(s_t, a_t))} d\tau_n \end{aligned}$$

(PAI, Levine, 2018)

- $p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})$ is the distribution of “optimal” trajectories!
 - Minimize KL-divergence between $p^{\pi_\theta}(\tau_n)$ and $p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})$

$$\begin{aligned}
 \mathcal{L}_\pi^{\text{PAI}} &:= D_{\text{KL}}[p^{\pi_\theta}(\tau_n) \parallel p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})] \\
 &\propto \int p^{\pi_\theta}(\tau_n) \ln \frac{\rho(s_0) \prod_{t=0}^{n-1} \pi_\theta(a_t | s_t) P(s_{t+1} | s_t, a_t)}{\rho(s_0) \prod_{t=0}^{n-1} \mu(a_t | s_t) P(s_{t+1} | s_t, a_t) \exp(\frac{1}{\alpha} r(s_t, a_t))} d\tau_n \\
 &= \int p^{\pi_\theta}(\tau_n) \ln \left(\prod_{t=0}^{n-1} \frac{\pi_\theta(a_t | s_t)}{\mu(a_t | s_t) \exp(\frac{1}{\alpha} r(s_t, a_t))} \right) d\tau_n
 \end{aligned}$$

(PAI, Levine, 2018)

- $p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})$ is the distribution of “optimal” trajectories!
 - Minimize KL-divergence between $p^{\pi_\theta}(\tau_n)$ and $p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})$

$$\begin{aligned}
 \mathcal{L}_\pi^{\text{PAI}} &:= D_{\text{KL}}[p^{\pi_\theta}(\tau_n) \parallel p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})] \\
 &\propto \int p^{\pi_\theta}(\tau_n) \ln \frac{\rho(s_0) \prod_{t=0}^{n-1} \pi_\theta(a_t | s_t) P(s_{t+1} | s_t, a_t)}{\rho(s_0) \prod_{t=0}^{n-1} \mu(a_t | s_t) P(s_{t+1} | s_t, a_t) \exp(\frac{1}{\alpha} r(s_t, a_t))} d\tau_n \\
 &= \int p^{\pi_\theta}(\tau_n) \ln \left(\prod_{t=0}^{n-1} \frac{\pi_\theta(a_t | s_t)}{\mu(a_t | s_t) \exp(\frac{1}{\alpha} r(s_t, a_t))} \right) d\tau_n \\
 &= - \int p^{\pi_\theta}(\tau_n) \sum_{t=0}^{n-1} \left(\frac{1}{\alpha} r(s_t, a_t) - \ln \frac{\pi_\theta(a_t | s_t)}{\mu(a_t | s_t)} \right) d\tau_n
 \end{aligned}$$

(PAI, Levine, 2018)

- $p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})$ is the distribution of “optimal” trajectories!
 - Minimize KL-divergence between $p^{\pi_\theta}(\tau_n)$ and $p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})$
 - using the uniform policy as prior μ yields entropy \mathcal{H}

$$\begin{aligned}
 \mathcal{L}_\pi^{\text{PAI}} &:= D_{\text{KL}}[p^{\pi_\theta}(\tau_n) \parallel p^\mu(\tau_n | \{O_t = 1\}_{t=0}^{n-1})] \\
 &\propto \int p^{\pi_\theta}(\tau_n) \ln \frac{\rho(s_0) \prod_{t=0}^{n-1} \pi_\theta(a_t | s_t) P(s_{t+1} | s_t, a_t)}{\rho(s_0) \prod_{t=0}^{n-1} \mu(a_t | s_t) P(s_{t+1} | s_t, a_t) \exp(\frac{1}{\alpha} r(s_t, a_t))} d\tau_n \\
 &= \int p^{\pi_\theta}(\tau_n) \ln \left(\prod_{t=0}^{n-1} \frac{\pi_\theta(a_t | s_t)}{\mu(a_t | s_t) \exp(\frac{1}{\alpha} r(s_t, a_t))} \right) d\tau_n \\
 &= - \int p^{\pi_\theta}(\tau_n) \sum_{t=0}^{n-1} \left(\frac{1}{\alpha} r(s_t, a_t) - \ln \frac{\pi_\theta(a_t | s_t)}{\mu(a_t | s_t)} \right) d\tau_n \\
 &\equiv - \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^{n-1} \left(r(s_t, a_t) + \alpha \mathcal{H}[\pi_\theta(\cdot | s_t)] \right) \right]
 \end{aligned}$$

- Stochastic policy gradient with additional entropy reward!

maximum entropy RL (PAI, [Levine, 2018](#)) is a lower bound to robust RL (as proven in [Eysenbach and Levine, 2021](#))

7.2 Stochastic vs. deterministic policy gradients



- Stochastic policy $\pi_\theta(a|s)$ from family $\hat{\pi}(a|\mathbf{f}_\theta(s))$
 - $\mathbf{f}_\theta(s) \in \mathbb{R}^m$ denotes the sufficient statistics at $s \in \mathcal{S}$
 - e.g. $\hat{\pi}(a|\mathbf{f}_\theta(s)) = \mathcal{N}(a | \mu_\theta(s), \sigma_\theta^2(s))$ with $\mathbf{f}_\theta(s) = [\mu_\theta(s), \sigma_\theta^2(s)]$

$$\nabla_\theta \mathcal{L}_\pi[\theta] \approx - \sum_{t=0}^{n-1} \gamma^t \iint \xi_t^\mu(s) \nabla_\theta \underbrace{\pi_\theta(a|s)}_{\hat{\pi}(a|\mathbf{f}_\theta(s))} Q^\pi(s, a) da ds$$

here we use $Q^\pi(s, a) = \mathbb{E}_\pi[R_0 | s_0 = s, a_0 = a]$, and the approximation drops $\frac{\xi_t^\pi(s_t)}{\xi_t^\mu(s_t)}$, see Sutton et al. (2016) for ways to estimate it

- Stochastic policy $\pi_\theta(a|s)$ from family $\hat{\pi}(a|\mathbf{f}_\theta(s))$
 - $\mathbf{f}_\theta(s) \in \mathbb{R}^m$ denotes the sufficient statistics at $s \in \mathcal{S}$
 - e.g. $\hat{\pi}(a|\mathbf{f}_\theta(s)) = \mathcal{N}(a | \mu_\theta(s), \sigma_\theta^2(s))$ with $\mathbf{f}_\theta(s) = [\mu_\theta(s), \sigma_\theta^2(s)]$

$$\begin{aligned}
 \nabla_\theta \mathcal{L}_\pi[\theta] &\approx - \sum_{t=0}^{n-1} \gamma^t \iint \xi_t^\mu(s) \underbrace{\nabla_\theta \pi_\theta(a|s)}_{\hat{\pi}(a|\mathbf{f}_\theta(s))} Q^\pi(s, a) da ds \\
 &= - \sum_{t=0}^{n-1} \gamma^t \int \xi_t^\mu(s) \underbrace{\nabla_\theta \mathbf{f}_\theta(s)}_{\boldsymbol{\pi}_\theta(s)} \underbrace{\nabla_{\mathbf{f}} \int \hat{\pi}(a|\mathbf{f}) Q^\pi(s, a) da}_{\hat{Q}^\pi(s, \mathbf{f})} \bigg|_{\mathbf{f}=\mathbf{f}_\theta(s)} ds
 \end{aligned}$$

Ciosek and Whiteson (2018) developed this view, called expected policy gradients (EPG)

7.2 Stochastic vs. deterministic policy gradients



- Stochastic policy $\pi_\theta(a|s)$ from family $\hat{\pi}(a|\mathbf{f}_\theta(s))$
 - $\mathbf{f}_\theta(s) \in \mathbb{R}^m$ denotes the sufficient statistics at $s \in \mathcal{S}$
 - e.g. $\hat{\pi}(a|\mathbf{f}_\theta(s)) = \mathcal{N}(a | \mu_\theta(s), \sigma_\theta^2(s))$ with $\mathbf{f}_\theta(s) = [\mu_\theta(s), \sigma_\theta^2(s)]$

$$\begin{aligned}
 \nabla_\theta \mathcal{L}_\pi[\theta] &\approx - \sum_{t=0}^{n-1} \gamma^t \iint \xi_t^\mu(s) \underbrace{\nabla_\theta \pi_\theta(a|s)}_{\hat{\pi}(a|\mathbf{f}_\theta(s))} Q^\pi(s, a) da ds \\
 &= - \sum_{t=0}^{n-1} \gamma^t \int \xi_t^\mu(s) \underbrace{\nabla_\theta \mathbf{f}_\theta(s)}_{\boldsymbol{\pi}_\theta(s)} \underbrace{\nabla_{\mathbf{f}} \int \hat{\pi}(a|\mathbf{f}) Q^\pi(s, a) da}_{\hat{Q}^\pi(s, \mathbf{f})} \bigg|_{\mathbf{f}=\mathbf{f}_\theta(s)} ds \\
 &= - \mathbb{E}_\mu \left[\sum_{t=0}^{n-1} \gamma^t \nabla_\theta \hat{Q}^\pi(s_t, \boldsymbol{\pi}_\theta(s_t)) \right]
 \end{aligned}$$

- SPG for family of distributions \equiv DPG on sufficient statistics!

Ciosek and Whiteson (2018) developed this view, called expected policy gradients (EPG)

7.2 The reparametrization trick



- Problem: estimate $\nabla_{\theta} \hat{Q}^{\pi}(s, \mathbf{f}_{\theta}(s)) = \int \nabla_{\theta} \hat{\pi}(\mathbf{a} | \mathbf{f}_{\theta}(s)) Q^{\pi}(s, \mathbf{a}) d\mathbf{a}$
 - hard to solve in general

this trick will be used in SQL on slide 19 and SAC on slide 20

7.2 The reparametrization trick



- Problem: estimate $\nabla_{\theta} \hat{Q}^{\pi}(s, \mathbf{f}_{\theta}(s)) = \int \nabla_{\theta} \hat{\pi}(\mathbf{a} | \mathbf{f}_{\theta}(s)) Q^{\pi}(s, \mathbf{a}) d\mathbf{a}$
 - hard to solve in general
- Bijective mappings can model arbitrary distributions
 - $\mathbf{a} := \mathbf{f}_{\theta}(s, \epsilon) \sim \pi_{\theta}(\cdot | s)$, $\epsilon \sim \mathcal{N}(\cdot | \mathbf{0}, \mathbf{I})$, for bijective \mathbf{f}_{θ}

$$\begin{aligned}\nabla_{\theta} \hat{Q}^{\pi}(s, \mathbf{f}_{\theta}(s)) &= \nabla_{\theta} \mathbb{E}[Q^{\pi}(s, \mathbf{a}) \mid \mathbf{a} \sim \pi_{\theta}(\cdot | s)] \\ &= \mathbb{E}[\nabla_{\theta} Q^{\pi}(s, \mathbf{f}_{\theta}(s, \epsilon)) \mid \epsilon \sim \mathcal{N}(\cdot | \mathbf{0}, \mathbf{I})] \\ &\approx \nabla_{\theta} \frac{1}{m} \sum_{i=1}^m Q^{\pi}(s, \mathbf{f}_{\theta}(s, \epsilon_i)), \quad \epsilon_i \sim \mathcal{N}(\cdot | \mathbf{0}, \mathbf{I})\end{aligned}$$

- One forward pass over batch with i.i.d. ϵ_i

this trick will be used in SQL on slide 19 and SAC on slide 20

- Robust RL chooses *disruption* that minimizes return
- Related planning as inference (PAI) is RL with entropy reward

⇒ Stochastic policies allow robust control!

- SPG \equiv DPG on sufficient statistics
- Reparametrization trick to compute SPG with DPG

Learning Objectives

- LO7.3: Explain robust RL and planning as inference
- LO7.4: Explain the connection between SPG and DPG
- LO7.5: Derive the reparametrization trick

7.3

Off-Policy Actor-Critic

Maximum entropy RL

7.3 Maximum entropy RL



- We can use additional entropy reward for policy gradient
- Or we can try to maximize the *soft Q-values* directly:

$$Q_{\text{soft}}^{\pi}(s, a) := r(s, a) + \mathbb{E}_{\pi} \left[\sum_{t=1}^{\infty} \gamma^t \left(r_t + \alpha \mathcal{H}[\pi(\cdot | s_t)] \right) \middle| \begin{matrix} s_0 = s \\ a_0 = a \end{matrix} \right]$$

$$Q_{\text{soft}}^*(s, a) := \max_{\pi} Q_{\text{soft}}^{\pi}(s, a) = r(s, a) + \gamma \mathbb{E} [V_{\text{soft}}^*(s') \mid s' \sim P(\cdot | s, a)]$$



7.3 Maximum entropy RL



- We can use additional entropy reward for policy gradient
- Or we can try to maximize the *soft Q-values* directly:

$$Q_{\text{soft}}^{\pi}(s, a) := r(s, a) + \mathbb{E}_{\pi} \left[\sum_{t=1}^{\infty} \gamma^t \left(r_t + \alpha \mathcal{H}[\pi(\cdot|s_t)] \right) \right] \Big|_{\substack{s_0=s \\ a_0=a}}$$

$$Q_{\text{soft}}^*(s, a) := \max_{\pi} Q_{\text{soft}}^{\pi}(s, a) = r(s, a) + \gamma \mathbb{E}[V_{\text{soft}}^*(s') \mid s' \sim P(\cdot|s, a)]$$

- There exist an analytical solution for V_{soft}^* and π_{soft}^* :

$$V_{\text{soft}}^*(s) := \alpha \ln \left(\int \exp \left(\frac{1}{\alpha} Q_{\text{soft}}^*(s, a) \right) da \right) \equiv \alpha \operatorname{softmax}_{a \in \mathcal{A}} \frac{1}{\alpha} Q_{\text{soft}}^*(s, a)$$

$$\pi_{\text{soft}}^*(a|s) := \arg \max_{\pi} Q_{\text{soft}}^{\pi}(s, a) = \exp \left(\frac{1}{\alpha} Q_{\text{soft}}^*(s, a) - \frac{1}{\alpha} V_{\text{soft}}^*(s) \right)$$

$$\max_a f(a) = \lim_{\alpha \rightarrow 0} \alpha \operatorname{softmax}_a \frac{1}{\alpha} f(a),$$

proof for V_{soft}^* and π_{soft}^* in [Ziebart \(2010\)](#) and/or [Haarnoja et al. \(2017, 2018a\)](#)

7.3 Soft Q-learning (SQL)



- One of the first attempts at maximum entropy RL
- SQL optimizes Q_ϕ and a corresponding policy π_θ
 - minimize KL-divergence between π_θ and π_{soft}^*
 - $V_\phi(s) := \alpha \ln \mathbb{E} \left[\frac{1}{\mu(a)} \exp \left(\frac{1}{\alpha} Q_\phi(s, a) \right) \middle| a \sim \mu \right]$

$$\mathcal{L}_Q^{\text{soft}}[\phi] := \mathbb{E} \left[\left(r + \gamma V_{\phi'}(s') - Q_\phi(s, a) \right)^2 \middle| \langle s, a, r, s' \rangle \sim \mathcal{D} \right]$$

$$\mathcal{L}_\pi^{\text{soft}}[\theta] := \mathbb{E} \left[D_{\text{KL}} \left[\pi_\theta(\cdot | s) \middle| \underbrace{\exp \left(\frac{1}{\alpha} Q_\phi(s, \cdot) - \frac{1}{\alpha} V_\phi(s) \right)}_{\pi_{\text{soft}}^*(\cdot | s)} \right] \middle| s \sim \mathcal{D} \right]$$

- SQL minimizes $\mathcal{L}_\pi^{\text{soft}}$ with *Stein variational gradient descent*
 - complicated and soon supplanted by SAC

(SQL, Haarnoja et al., 2017)

- Same basic losses as SQL, but with a ton of tricks
 - extra state-value function $V_\psi(s)$
 - TD3 style twin Q-values $\bar{Q}(s, a) := \min_{i \in \{1, 2\}} Q_{\phi_i}(s, a)$

$$\mathcal{L}_V^{\text{soft}}[\psi] := \mathbb{E} \left[\left(\mathbb{E}[\bar{Q}(s, a) - \alpha \ln \pi_\theta(a|s) \mid a \sim \pi_\theta(\cdot|s)] - V_\psi(s) \right)^2 \mid s \sim \mathcal{D} \right]$$

$$\mathcal{L}_Q^{\text{soft}}[\phi] := \mathbb{E} \left[\left(r + \gamma V_{\psi'}(s') - Q_\phi(s, a) \right)^2 \mid \langle s, a, r, s' \rangle \sim \mathcal{D} \right]$$

$$\mathcal{L}_\pi^{\text{soft}}[\theta] := \mathbb{E} \left[D_{\text{KL}}[\pi_\theta(\cdot|s) \parallel \exp \left(\frac{1}{\alpha} \bar{Q}(s, \cdot) - \frac{1}{\alpha} V_\psi(s) \right)] \mid s \sim \mathcal{D} \right]$$

- SAC implementations branch off TD3!
 - uses replay buffer, pessimism, but not clipped noise
 - stochastic policy via reparametrization trick (slide 16)
 - action squashing for bounded action spaces
 - automatic entropy adjustment

7.3 Action squashing



- Gaussian distribution not suitable for bounded action spaces
 - use Gaussian $f_{\theta}(s, \epsilon) := \mu_{\theta}(s) + \sigma_{\theta}^2(s) \odot \epsilon$, $\epsilon \sim \mathcal{N}(\cdot | \mathbf{0}, \mathbf{I})$
 - “squash” action into $(-1, 1)$: $\mathbf{a} := \tanh(f_{\theta}(s, \epsilon))$
- How does this change $\ln \pi(\mathbf{a} | s)$ for the entropy?

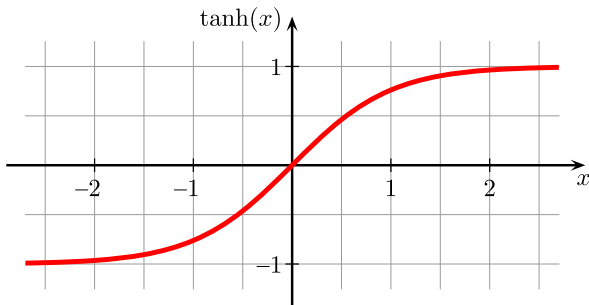


image source [wikipedia.org](https://en.wikipedia.org/wiki/Hyperbolic_tangent)

7.3 Action squashing



- Gaussian distribution not suitable for bounded action spaces
 - use Gaussian $\mathbf{f}_\theta(s, \epsilon) := \boldsymbol{\mu}_\theta(s) + \boldsymbol{\sigma}_\theta^2(s) \odot \epsilon$, $\epsilon \sim \mathcal{N}(\cdot | \mathbf{0}, \mathbf{I})$
 - “squash” action into $(-1, 1)$: $\mathbf{a} := \tanh(\mathbf{f}_\theta(s, \epsilon))$
- How does this change $\ln \pi(\mathbf{a}|s)$ for the entropy?
 - $\pi(\mathbf{a}|s) = \mathcal{N}\left(\tanh^{-1}(\mathbf{a}) \mid \boldsymbol{\mu}_\theta(s), \boldsymbol{\sigma}_\theta^2(s)\right) \left| \det \left[\frac{\partial \tanh^{-1}(\mathbf{a}')}{\partial \mathbf{a}'} \Big|_{\mathbf{a}'=\mathbf{a}} \right] \right|$

https://en.wikipedia.org/wiki/Probability_density_function#Vector_to_vector

7.3 Action squashing



- Gaussian distribution not suitable for bounded action spaces

- use Gaussian $\mathbf{f}_\theta(s, \epsilon) := \boldsymbol{\mu}_\theta(s) + \boldsymbol{\sigma}_\theta^2(s) \odot \epsilon$, $\epsilon \sim \mathcal{N}(\cdot | \mathbf{0}, \mathbf{I})$
- “squash” action into $(-1, 1)$: $\mathbf{a} := \tanh(\mathbf{f}_\theta(s, \epsilon))$

- How does this change $\ln \pi(\mathbf{a}|s)$ for the entropy?

- $\pi(\mathbf{a}|s) = \mathcal{N}\left(\tanh^{-1}(\mathbf{a}) \mid \boldsymbol{\mu}_\theta(s), \boldsymbol{\sigma}_\theta^2(s)\right) \left| \det \left[\frac{\partial \tanh^{-1}(\mathbf{a}')}{\partial \mathbf{a}'} \Big|_{\mathbf{a}'=\mathbf{a}} \right] \right|$
- inverse function th.: $\frac{\partial \tanh^{-1}(\mathbf{a}')}{\partial \mathbf{a}'} \Big|_{\mathbf{a}'=\mathbf{a}} = \left(\frac{\partial \tanh(\mathbf{u})}{\partial \mathbf{u}} \Big|_{\mathbf{u}=\tanh^{-1}(\mathbf{a})} \right)^{-1}$

https://en.wikipedia.org/wiki/Inverse_function_theorem

7.3 Action squashing



- Gaussian distribution not suitable for bounded action spaces

- use Gaussian $\mathbf{f}_\theta(s, \epsilon) := \boldsymbol{\mu}_\theta(s) + \boldsymbol{\sigma}_\theta^2(s) \odot \epsilon$, $\epsilon \sim \mathcal{N}(\cdot | \mathbf{0}, \mathbf{I})$
- “squash” action into $(-1, 1)$: $\mathbf{a} := \tanh(\mathbf{f}_\theta(s, \epsilon))$

- How does this change $\ln \pi(\mathbf{a}|s)$ for the entropy?

- $\pi(\mathbf{a}|s) = \mathcal{N}\left(\tanh^{-1}(\mathbf{a}) \mid \boldsymbol{\mu}_\theta(s), \boldsymbol{\sigma}_\theta^2(s)\right) \left| \det \left[\frac{\partial \tanh^{-1}(\mathbf{a}')}{\partial \mathbf{a}'} \Big|_{\mathbf{a}'=\mathbf{a}} \right] \right|$
- inverse function th.: $\frac{\partial \tanh^{-1}(\mathbf{a}')}{\partial \mathbf{a}'} \Big|_{\mathbf{a}'=\mathbf{a}} = \left(\frac{\partial \tanh(\mathbf{u})}{\partial \mathbf{u}} \Big|_{\mathbf{u}=\tanh^{-1}(\mathbf{a})} \right)^{-1}$
- determinant of inverse matrix: $\det \mathbf{J}^{-1} = (\det \mathbf{J})^{-1}$

https://en.wikipedia.org/wiki/Determinant#Multiplicativity_and_matrix_groups

7.3 Action squashing



- Gaussian distribution not suitable for bounded action spaces

- use Gaussian $f_{\theta}(s, \epsilon) := \mu_{\theta}(s) + \sigma_{\theta}^2(s) \odot \epsilon$, $\epsilon \sim \mathcal{N}(\cdot | \mathbf{0}, \mathbf{I})$
- “squash” action into $(-1, 1)$: $\mathbf{a} := \tanh(f_{\theta}(s, \epsilon))$

- How does this change $\ln \pi(\mathbf{a}|s)$ for the entropy?

- $\pi(\mathbf{a}|s) = \mathcal{N}\left(\tanh^{-1}(\mathbf{a}) \mid \mu_{\theta}(s), \sigma_{\theta}^2(s)\right) \left| \det \left[\frac{\partial \tanh^{-1}(\mathbf{a}')}{\partial \mathbf{a}'} \Big|_{\mathbf{a}'=\mathbf{a}} \right] \right|$
- inverse function th.: $\frac{\partial \tanh^{-1}(\mathbf{a}')}{\partial \mathbf{a}'} \Big|_{\mathbf{a}'=\mathbf{a}} = \left(\frac{\partial \tanh(\mathbf{u})}{\partial \mathbf{u}} \Big|_{\mathbf{u}=\tanh^{-1}(\mathbf{a})} \right)^{-1}$
- determinant of inverse matrix: $\det \mathbf{J}^{-1} = (\det \mathbf{J})^{-1}$
- diagonal Jacobian matrix $J_{kk} := \frac{\partial \tanh(u_k)}{\partial u_k} = (1 - \tanh^2(u_k))$
 $\Rightarrow \det \mathbf{J} = \prod_k J_{kk}$

https://en.wikipedia.org/wiki/Determinant#Immediate_consequences

7.3 Action squashing



- Gaussian distribution not suitable for bounded action spaces

- use Gaussian $\mathbf{f}_\theta(s, \epsilon) := \boldsymbol{\mu}_\theta(s) + \boldsymbol{\sigma}_\theta^2(s) \odot \epsilon$, $\epsilon \sim \mathcal{N}(\cdot | \mathbf{0}, \mathbf{I})$
- “squash” action into $(-1, 1)$: $\mathbf{a} := \tanh(\mathbf{f}_\theta(s, \epsilon))$

- How does this change $\ln \pi(\mathbf{a}|s)$ for the entropy?

- $\pi(\mathbf{a}|s) = \mathcal{N}\left(\tanh^{-1}(\mathbf{a}) \mid \boldsymbol{\mu}_\theta(s), \boldsymbol{\sigma}_\theta^2(s)\right) \left| \det \left[\frac{\partial \tanh^{-1}(\mathbf{a}')}{\partial \mathbf{a}'} \Big|_{\mathbf{a}'=\mathbf{a}} \right] \right|$
- inverse function th.: $\frac{\partial \tanh^{-1}(\mathbf{a}')}{\partial \mathbf{a}'} \Big|_{\mathbf{a}'=\mathbf{a}} = \left(\frac{\partial \tanh(\mathbf{u})}{\partial \mathbf{u}} \Big|_{\mathbf{u}=\tanh^{-1}(\mathbf{a})} \right)^{-1}$
- determinant of inverse matrix: $\det \mathbf{J}^{-1} = (\det \mathbf{J})^{-1}$
- diagonal Jacobian matrix $J_{kk} := \frac{\partial \tanh(u_k)}{\partial u_k} = (1 - \tanh^2(u_k))$
 $\Rightarrow \det \mathbf{J} = \prod_k J_{kk}$

- $\ln \pi(\mathbf{a}|s) = \underbrace{\ln \mathcal{N}(\tanh^{-1}(\mathbf{a}) \mid \boldsymbol{\mu}_\theta(s), \boldsymbol{\sigma}_\theta^2(s))}_{\text{Gaussian entropy}} - \underbrace{\sum_k \ln(1 - \tanh^2(a_k))}_{\text{correction term for squashing}}$

action squashing can be numerically unstable

see also [Haarnoja et al. \(2018b, Appendix C\)](#)

7.3 Automatic entropy adjustment



- Entropy hyper-parameter α is hard to choose
 - too low and the entropy term is insignificant
 - too high and the reward is insignificant

- Automatic adjustment by defining target entropy H

$$\min_{\theta} \mathcal{L}_{\pi}[\theta] \quad \text{s.t.} \quad -\mathbb{E} \left[\ln \pi_{\theta}(a|s) \mid a \sim \pi_{\theta}^{\mathcal{D}}(\cdot|s) \right] \geq H$$

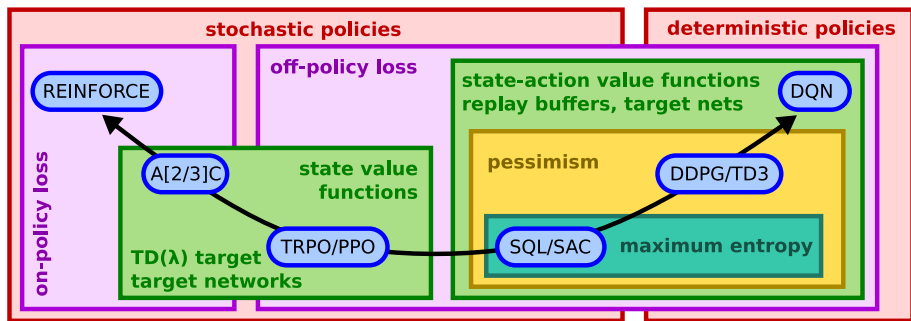
- Gradient descend w.r.t. dual variable α of the Lagrangian
 - averaged over minibatch $\{s_t\}_{i=1}^n$ and $a_t^i \sim \pi_{\theta}(\cdot|s_t), \forall i, t$
 - $\alpha \leftarrow \alpha + \eta \mathbb{E} \left[\ln \pi_{\theta}(a|s) + H \mid a \sim \pi_{\theta}^{\mathcal{D}}(\cdot|s) \right] \approx \alpha + \frac{\eta}{nm} \sum_{t=1}^n \sum_{i=1}^m \ln \pi_{\theta}(a_t^i|s_t) + \eta H$

based on derivation in [Haarnoja et al. \(2018b\)](#)

7.3 SQL/SAC in comparison



- Entropy maximization for robustness and exploration
 - reparameterization trick to extend TD3 to stochastic policies
 - state-action Q-value functions generalize worse (like TD3)
 - SAC and TD3 stabilize Q-value with pessimism



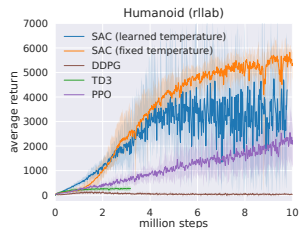
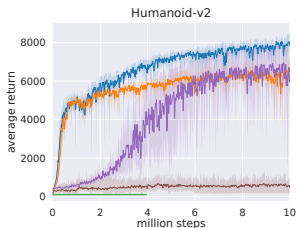
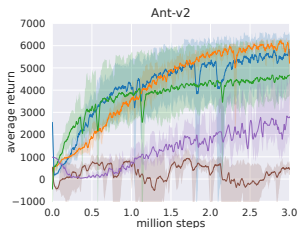
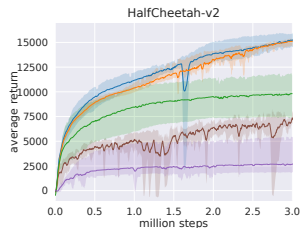
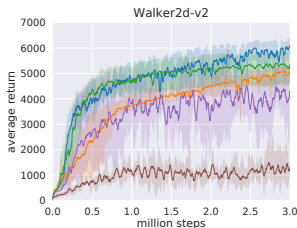
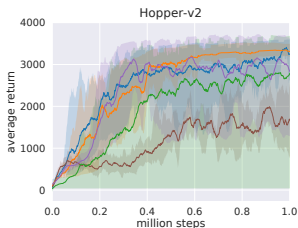
SAC is also one of the best RL algorithms, explores slightly better than TD3 and is very common in *simulated* robotics

7.3 Performance comparison



- AC2/3 < DDPG < PPO $\stackrel{?}{\leq}$ (TD3 $\stackrel{?}{\approx}$ SAC)

(no consensus)



results from SAC paper ([Haarnoja et al., 2018b](#)), for results that favor TD3, see slide 8

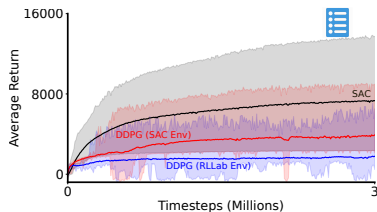
7.3 Policy improvement in DRL algorithms



- Pure on-policy algorithms (slow, stable, require state-value v_ϕ)
 - Reinforce $\mathcal{L}_\pi^{\text{RE}}[\theta] = -\frac{1}{n} \sum_{t=1}^n \gamma^t R_t \ln \pi_\theta(a_t|s_t)$ (high variance)
 - A(|2|3)C $\mathcal{L}_\pi^{\text{AC}}[\theta] = -\frac{1}{n} \sum_{t=1}^n \gamma^t \underbrace{(y_t^{n/\lambda}(\tau_n) - V_\phi(s_t))}_{A_t} \ln \pi_\theta(a_t|s_t)$
- On-off-policy algorithms (fast, need trust-region for stability)
 - Off-PAC $\mathcal{L}_\mu^{\text{OPAC}}[\theta] = -\frac{1}{n} \sum_{t=1}^n \gamma^t A_t \frac{\pi_\theta(a_t|s_t)}{\mu(a_t|s_t)}$ (extremely unstable)
 - TRPO $\mathcal{L}_\mu^{\text{TRPO}}[\theta] = -\frac{1}{n} \sum_{t=1}^n \gamma^t A_t \frac{\pi_\theta(a_t|s_t)}{\mu(a_t|s_t)}$ s.t. $\frac{1}{n} \sum_{t=1}^n D_{\text{KL}}(\mu(\cdot|s_t) \parallel \pi_\theta(\cdot|s_t))$
 - PPO $\mathcal{L}_\mu^{\text{PPO}}[\theta] = -\frac{1}{n} \sum_{t=1}^n \gamma^t \min\left(A_t \frac{\pi_\theta(a_t|s_t)}{\mu(a_t|s_t)}, A_t \text{clip}\left(\frac{\pi_\theta(a_t|s_t)}{\mu(a_t|s_t)}, 1-\epsilon, 1+\epsilon\right)\right)$
- Off-policy algorithms (fast, require Q -value, easily overestimated)
 - DDPG/TD3 $\mathcal{L}_\mu^{\text{DDPG}}[\theta] = -\frac{1}{n} \sum_{t=1}^n \gamma^t \hat{Q}_\phi^\pi(s_t, \pi(s_t))$ (Q -value of suff. statistics)
 - SAC $\mathcal{L}_\pi^{\text{soft}}[\theta] = \frac{1}{n} \sum_{t=1}^n D_{\text{KL}}\left[\pi_\theta(\cdot|s) \parallel \exp\left(\frac{1}{\alpha} \min_i Q_{\phi_i}(s_t, \cdot) - \frac{1}{\alpha} V_\psi(s_t)\right)\right]$
 - DQN $\pi(a|s_t) = \arg \max_a Q_\phi^*(s_t, a)$ (greedy Q -values Q^*)

7.3 A word of caution

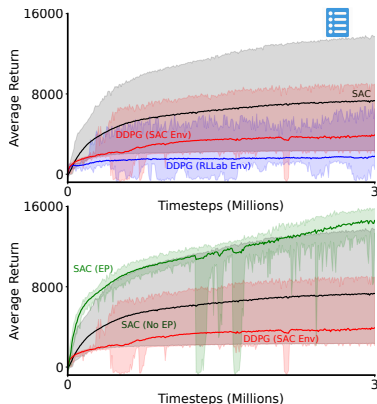
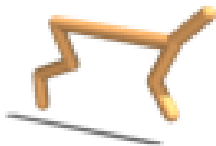
- Replicating SAC results ([Haarnoja et al., 2018a](#))
 - Mujoco's Half-Cheetha env
 - very unreliable training [30 seeds]



figures and story are from [Patterson et al. \(2023\)](#), which is a good take on best practices in DRL

7.3 A word of caution

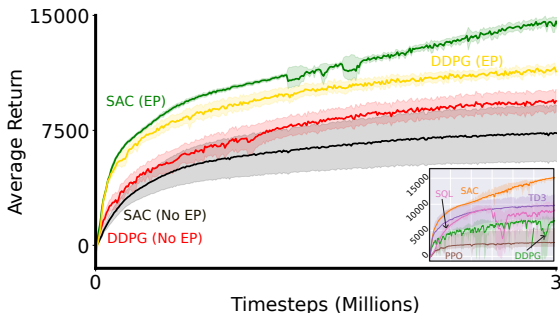
- Replicating SAC results ([Haarnoja et al., 2018a](#))
- Found undocumented feature
 - random initial “burn-in” phase (EP)
 - performance now reliable



figures and story are from [Patterson et al. \(2023\)](#), which is a good take on best practices in DRL

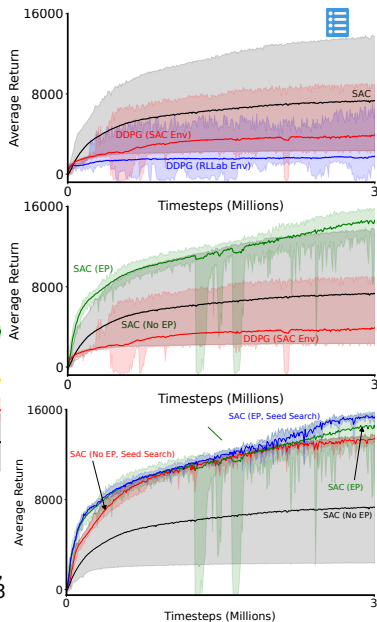
7.3 A word of caution

- Replicating SAC results (Haarnoja et al., 2018a)
- Found undocumented feature
- Cherry-picking seeds (**don't do this!**)
 - same performance as (Haarnoja et al., 2018a)



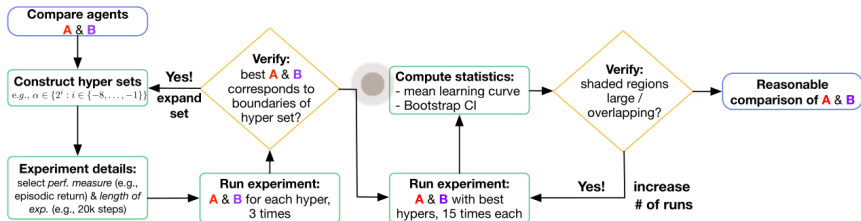
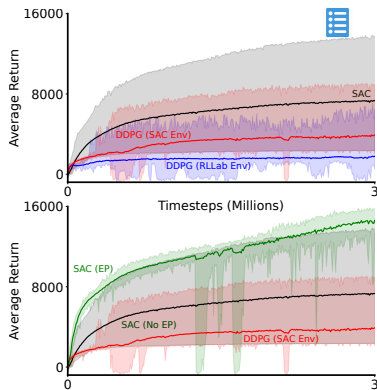
Timesteps (Millions)

figures and story are from [Patterson et al. \(2023\)](#), which is a good take on best practices in DRL



7.3 A word of caution

- Replicating SAC results (Haarnoja et al., 2018a)
- Found undocumented feature
- Cherry-picking seeds (**don't do this!**)
- Patterson et al. (2023) recommends:



figures and story are from [Patterson et al. \(2023\)](#), which is a good take on best practices in DRL

- Maximum entropy RL yields a softmax solution
- SQL implements this with KL-divergence
- SAC improves by branching off TD3 with reparametrization
- SAC also uses action squashing and entropy adjustment
- Empirical comparison inconclusive, TD3 and SAC both SOTA

Learning Objectives

LO7.6: Explain maximum entropy RL, SQL and SAC

LO7.7: Explain action squashing and automatic entropy adjustment

LO7.8: Discuss commonalities and differences between policy gradients

7.3 Next lecture



- Next lecture: **exploration!**
- Start with assignment sheet 3
- Questions? Ask them here:
answers.ewi.tudelft.nl



SO MUCH OF "AI" IS JUST FIGURING OUT WAYS TO OFFLOAD WORK ONTO RANDOM STRANGERS.

image source: xkcd.com

- Kamil Ciosek and Shimon Whiteson. Expected policy gradients. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, pages 2868–2875, 2018. URL <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16116>.
- Benjamin Eysenbach and Sergey Levine. Maximum entropy rl (provably) solves some robust rl problems, 2021. URL <https://arxiv.org/abs/2103.06257>.
- Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *Proceedings of Machine Learning Research (ICML)*, volume 80, pages 1587–1596, 2018. URL <http://proceedings.mlr.press/v80/fujimoto18a.html>.
- Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with deep energy-based policies. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, 2017. URL <https://proceedings.mlr.press/v70/haarnoja17a>.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of Machine Learning Research (ICML)*, volume 80, pages 1861–1870, 2018a. URL <http://proceedings.mlr.press/v80/haarnoja18b.html>.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Soft actor-critic algorithms and applications. *CoRR*, abs/1812.05905, 2018b. URL <http://arxiv.org/abs/1812.05905>.
- Sergey Levine. Reinforcement learning and control as probabilistic inference: Tutorial and review. *CoRR*, abs/1805.00909, 2018. URL <https://arxiv.org/abs/1805.00909>.
- Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2016. URL <http://arxiv.org/abs/1509.02971>.
- Zhuang Liu, Xuanlin Li, Bingyi Kang, and Trevor Darrell. Regularization matters in policy optimization - an empirical study on continuous control. In *International Conference on Learning Representations*, 2021. URL <https://arxiv.org/abs/1910.09191>.



- Janosch Moos, Kay Hansel, Hany Abdulsamad, Svenja Stark, Debora Clever, and Jan Peters. Robust reinforcement learning: A review of foundations and recent advances. *Machine Learning and Knowledge Extraction*, 4(1):276–315, 2022. ISSN 2504-4990. URL <https://doi.org/10.3390/make4010013>.
- Jun Morimoto and Kenji Doya. Robust reinforcement learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 13. MIT Press, 2001. URL <https://proceedings.neurips.cc/paper/2000/hash/e8dfff4676a47048d6f0c4ef899593dd-Abstract.html>.
- Andrew Patterson, Samuel Neumann, Martha White, and Adam White. Empirical design in reinforcement learning, 2023. URL <https://arxiv.org/abs/2304.01315>.
- David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy gradient algorithms. In *Proceedings of Machine Learning Research (ICML)*, volume 32, pages 387–395, 2014. URL <http://proceedings.mlr.press/v32/silver14.html>.
- Richard S. Sutton, A. Rupam Mahmood, and Martha White. An emphatic approach to the problem of off-policy temporal-difference learning. *Journal of Machine Learning Research*, 17(1):2603–2631, jan 2016. ISSN 1532-4435. URL <https://arxiv.org/abs/1503.04269>.
- Brian D. Ziebart. *Modeling Purposeful Adaptive Behavior with the Principle of Maximum Causal Entropy*. PhD thesis, Carnegie Mellon University, USA, 2010. URL <https://www.cs.cmu.edu/~bziebart/publications/thesis-bziebart.pdf>.