

Theory and Methods for Reinforcement Learning

Prof. Volkan Cevher
volkan.cevher@epfl.ch

Lecture 8: Deep and Robust Reinforcement Learning

Laboratory for Information and Inference Systems (LIONS)
École Polytechnique Fédérale de Lausanne (EPFL)

EE-618 (Spring 2023)



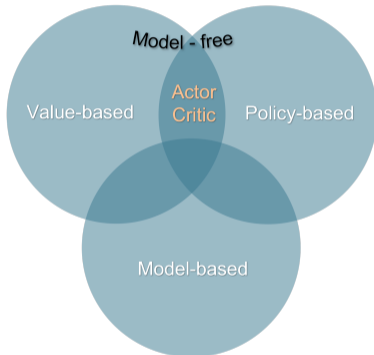
License Information for Theory and Methods for Reinforcement Learning (EE-618)

- ▷ This work is released under a [Creative Commons License](#) with the following terms:
- ▷ **Attribution**
 - ▶ The licensor permits others to copy, distribute, display, and perform the work. In return, licensees must give the original authors credit.
- ▷ **Non-Commercial**
 - ▶ The licensor permits others to copy, distribute, display, and perform the work. In return, licensees may not use the work for commercial purposes – unless they get the licensor's permission.
- ▷ **Share Alike**
 - ▶ The licensor permits others to distribute derivative works only under a license identical to the one that governs the licensor's work.
- ▷ [Full Text of the License](#)

Recap: Overview of reinforcement learning approaches

Value-based RL (Critic-only)

- Learn the optimal value functions V^* , Q^*
- **Algorithms:** Monte Carlo, SARSA, Q-learning, etc.
- Use temporal difference (low variance)
- Does not scale to large action spaces

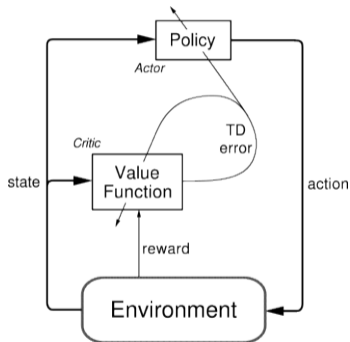


Policy-based RL (Actor-only)

- Learn the optimal policy via gradient methods
- **Algorithms:** PG, NPG, TRPO, PPO, etc.
- Scales to large or continuous action spaces
- High variance, sample inefficiency

Actor-Critic (AC) methods

- AC methods aim at combining the advantages of actor-only methods and critic-only methods.



Interaction of Actor-Critic [25].

- The actor uses the policy gradient to update the learning policy.
- The critic uses temporal difference learning to estimate the value function.

Actor-Critic methods

- Actor-critic algorithms follow an approximate policy gradient:

$$\nabla_{\theta} J(\pi_{\theta}) \approx \frac{1}{1-\gamma} \mathbb{E}_{s \sim \lambda_{\mu}^{\pi_{\theta}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [Q_w(s, a) \nabla_{\theta} \log \pi_{\theta}(a | s)] \right].$$

$$\nabla_{\theta} J(\pi_{\theta}) \approx \frac{1}{1-\gamma} \mathbb{E}_{s \sim \lambda_{\mu}^{\pi_{\theta}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [A_w(s, a) \nabla_{\theta} \log \pi_{\theta}(a | s)] \right].$$

- Actor: adjust the policy parameter θ using policy gradient using the value function estimated by the critic.
- Critic: update the parameter w to estimate action-value or advantage function.

$$Q_w(s, a) \approx Q^{\pi_{\theta}}(s, a)$$

$$A_w(s, a) \approx Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s)$$

Bias in Actor-Critic methods

- Recall action value expression of policy gradient

$$\nabla_{\theta} J(\pi_{\theta}) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim \lambda_{\mu}^{\pi_{\theta}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)} [Q^{\pi_{\theta}}(s, a) \nabla_{\theta} \log \pi_{\theta}(a | s)] \right].$$

- Policy gradient estimators used by actor-critic algorithms:

$$\hat{\nabla}_{\theta} J(\pi_{\theta}) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim \lambda_{\mu}^{\pi_{\theta}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)} [Q_w(s, a) \nabla_{\theta} \log \pi_{\theta}(a | s)] \right].$$

- Approximating the policy gradient using value function approximation Q_w could introduce bias.
- Luckily, if the value function approximation Q_w is chosen carefully, one may avoid such bias.

Compatible function approximation theorem

Compatible function approximation theorem [26]

Suppose the following two conditions are satisfied:

- Value function approximation at w^* is compatible to the policy, i.e.,

$$\nabla_w Q_{w^*}(s, a) = \nabla_\theta \log \pi_\theta(a | s).$$

- Value function parameter w^* minimizes the mean-squared error, i.e.,

$$\min_w \mathbb{E}_{s \sim \lambda_\mu^{\pi_\theta}, a \sim \pi_\theta(\cdot | s)} [(Q_w(s, a) - Q^{\pi_\theta}(s, a))^2].$$

Then the policy gradient using critic $Q_{w^*}(s, a)$ is exact:

$$\nabla_\theta J(\theta) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim \lambda_\mu^{\pi_\theta}, a \sim \pi_\theta(\cdot | s)} [\nabla_\theta \log \pi_\theta(a | s) Q_{w^*}(s, a)].$$

Remarks:

- Proof follows immediately from first-order optimality condition.
- Example: $Q_w(s, a) = \nabla_\theta \log \pi_\theta(a | s)^\top w$.

Variant I: Online Action-Value Actor-Critic

Online Action-Value Actor-Critic Algorithm

Initialize θ_0, w_0 , state $s_0 \sim \mu, a_0 \sim \pi_{\theta_0}(\cdot | s_0)$.

for each step of the episode $t = 0, \dots, T$ **do**

Obtain (r_t, s_{t+1}, a_{t+1}) from π_{θ_t} .

Compute policy gradient estimator: $\hat{\nabla}_{\theta} J(\pi_{\theta_t}) = Q_{w_t}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta_t}(a_t | s_t)$.

Actor update θ : $\theta_{t+1} = \theta_t + \alpha_t \hat{\nabla}_{\theta} J(\pi_{\theta_t})$.

Compute temporal difference: $\delta_t = r_t + \gamma Q_{w_t}(s_{t+1}, a_{t+1}) - Q_{w_t}(s_t, a_t)$.

Critic update: $w_{t+1} = w_t - \beta_t \delta_t \nabla_w Q_{w_t}(s_t, a_t)$.

end for

Remarks:

- Uses temporal difference to estimate the value function $Q^{\pi_{\theta}}$.
- Examples for Q_w : linear value function approximation $Q_w(s, a) = \phi(s, a)^{\top} w$.

Variante II: Advantage Actor-Critic

Advantage Actor-Critic (A2C)

Initialize θ_0 , w_0 , state $s_0 \sim \mu$.

for each step of the episode $t = 0, \dots, T$ **do**

Take action $a_t \sim \pi_{\theta_t}(\cdot | s_t)$, obtain (r_t, s_{t+1}) .

Estimate advantage function: $\delta_t = r_t + \gamma V_{w_t}(s_{t+1}) - V_{w_t}(s_t)$.

Compute policy gradient estimator: $\hat{\nabla}_{\theta} J(\pi_{\theta_t}) = \delta_t \nabla_{\theta} \log \pi_{\theta_t}(a_t | s_t)$.

Actor update: $\theta_{t+1} = \theta_t + \alpha_t \hat{\nabla}_{\theta} J(\pi_{\theta_t})$.

Critic update: $w_{t+1} = w_t - \beta_t \delta_t \nabla_w V_{w_t}(s_t)$.

end for

Remarks:

- Use $V_w(s)$ to approximate $V^{\pi_{\theta}}(s)$, for instance $V^w(s) \approx \phi(s)^{\top} w$.
- Use one step lookahead to estimate $Q^{\pi_{\theta}}(s_t, a_t) \approx r(s_t, a_t) + \gamma V^{\pi_{\theta}}(s_{t+1})$.
- Use advantage function to approximate the policy gradient.

Various Actor-Critic extensions

- **Natural Actor-Critic [17]**: use TRPO[22] or NPG[9] to update the actor
- **Actor-Critic with generalized advantage estimator [23]**: generalize advantage function with TD(λ)

$$\hat{A}_t^k(s_t, a_t) = r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \dots + \gamma^k V_w(s_{t+k}) - V_w(s_t)$$

$$\hat{A}_t^{\text{GAE}}(s_t, a_t) = (1 - \lambda) \sum_{k=1}^{\infty} \lambda^{k-1} \hat{A}_t^k(s_t, a_t)$$

- **Soft Actor-Critic [7]**: use entropy regularization in the objective to improve exploration

$$\max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) + \lambda \cdot \mathcal{H}(\pi(\cdot|s_t)) \right], \text{ where } \mathcal{H}(\pi(\cdot|s)) = \mathbb{E}_{a \sim \pi(\cdot|s)} [-\log \pi(a|s)]$$

Convergence of Actor-Critic methods

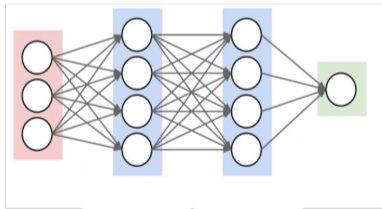
Remarks:

- The asymptotic analysis of two time-scale actor-critic methods (i.e., $\lim_{t \rightarrow \infty} \frac{\alpha_t}{\beta_t} = 0$) was established in [3] and [11].
- The proof is based on two-time-scale stochastic approximation and ODE analysis.
- Finite-sample analyses of actor-critic methods (tabular or LFA) have been studied very recently.
- This work is based on the bilevel optimization perspective; see e.g., [34].
- Indeed, Actor-critic algorithms can be formulated as bilevel optimization:

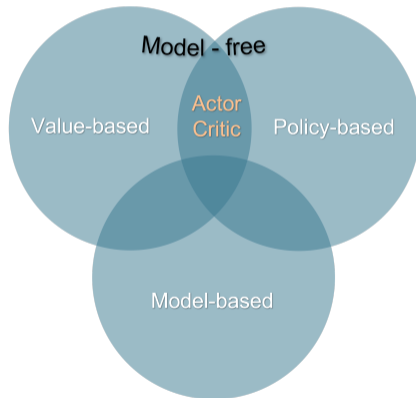
$$\begin{aligned} \min_{\theta} \quad & F(\theta) = f(\theta, w^*(\theta)), \quad (\text{Upper level}) \\ \text{s.t.} \quad & w^*(\theta) \in \arg \min_w \ell(\theta, w). \quad (\text{Lower level}) \end{aligned}$$

Deep reinforcement learning = DL + RL

- Tabular methods and linear function approximation are insufficient for large-scale RL applications.
- Using neural networks seems to be a must.

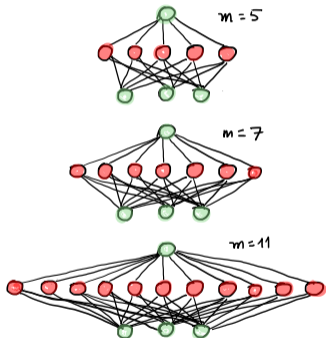


+



Neural networks

- Nested composition of (learnable) linear transformation with (fixed) nonlinear activation functions
- Example: a single-layer neural network (shallow neural network)



$$f(\mathbf{x}; W, \alpha) = \sum_{i=1}^m \alpha_i \cdot \sigma(w_i^\top \mathbf{x})$$

Activation function $\sigma(\cdot)$

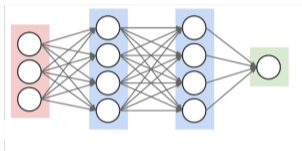
- Identity: $\sigma(u) = u$
- Sigmoid: $\sigma(u) = \frac{1}{1 + \exp(-u)}$
- Tanh: $\sigma(u) = \tanh(u)$
- Rectified linear unit (ReLU): $\sigma(u) = \max(0, u)$
-

Figure: Networks of increasing width

Deep neural networks

- More hidden layers, different activation functions, more general graph structure

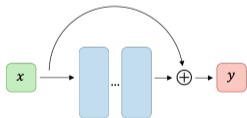
Feed forward network



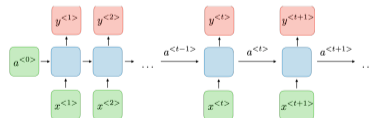
Convolutional network



Residual network



Recurrent network



Why neural networks?

- **Universal Approximation**

- Any continuous function on a compact domain can be (uniformly) approximated to arbitrary accuracy by a single-hidden layer neural network with a non-polynomial activation function. [Cybenko, 1989; Hornik et al., 1989; Barron, 1993]
- But the number of neurons can be large.

- **Benefits of depth**

- A deep network cannot be approximated by a reasonably-sized shallow network.[35]
- For example, there exists a function with $O(L^2)$ layers and width 2 which requires width $O(2^L)$ to approximate with $O(L)$ layers [27]. For more refined depth separation results see [20].

Example: ATARI network architecture

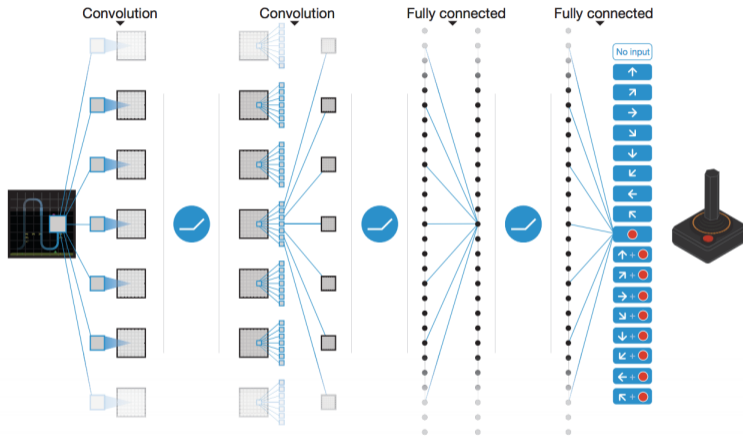


Figure: ATARI Network Architecture for $Q(s, a)$: History of frames as input. One output per action. [14]

Challenges with training neural networks in RL

- Deadly triad (Divergence when combining function approximation, bootstrapping, and off-policy learning)
- Non i.i.d. data
- Sample inefficiency
- High variance
- Overfitting
- Saddle points
- ...

Common Fixes or RL Tricks

- **Better data:** e.g., experience replay (mix online data and a buffer from past experience)
 - Reduce correlation, allow mini-batch update
- **Better objective:** e.g., use entropy regularization
 - Improve optimization landscape, encourage exploration
- **Better optimizers:** e.g., adaptive SGD such as Adam and RMSProp
 - Adaptive learning rates
- **Better estimation:** e.g., use eligibility traces, target works
 - Reduce overestimation bias, balance bias-variance tradeoff
- **Better sampling:** e.g., use prioritized replay (sample based on priority)
 - Prioritize transitions on which we can learn much
- **Better implementation:** e.g., parallel implementation (multithreading of CPU)
 - Speed up training, reduce correlation, allow better exploration
- **Better architectures:** e.g. dueling networks
 - Encode inductive biases that are good for RL

Value-based DRL

- Idea: use neural networks for value function approximation
- Recall Q-learning:

Q Learning

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha_t [r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

Q-learning with function approximation

$$w_{t+1} \leftarrow w_t + \alpha_t [r_t + \gamma \max_a Q_{w_t}(s_{t+1}, a) - Q_{w_t}(s_t, a_t)] \nabla Q_{w_t}(s_t, a_t)$$

- Note that Q-learning is not a stochastic gradient descent method.
- Naive deep Q-learning could diverge due to sample correlation and moving targets.
- **Deep Q-Networks (DeepMind, 2015) [14]**: combine several techniques for stabilizing Q-learning
 - Experience replay (better data efficiency and make data more stationary)
 - Target networks (prevent target objective from changing too fast)

Deep Q-Networks (DQN)

- Main idea: minimize the following mean-square error by SGD (or adaptive SGD)

$$\min_w L(w) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}} \left[\left(r + \gamma \max_{a'} Q(s', a'; w^-) - Q(s, a; w) \right)^2 \right]$$

- The target parameter w^- is held fixed and updated periodically

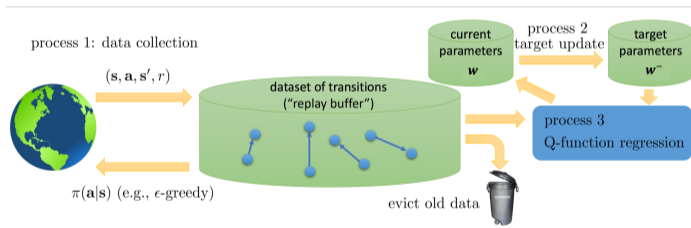


Figure: A more general view of DQN. Source: <https://zhuanlan.zhihu.com/p/468385820>

DQN in playing Atari games [14]



Figure: Five Atari 2600 Games: Pong, Breakout, Space Invaders, Seaquest, Beam Rider

	B. Rider	Breakout	Enduro	Pong	Q*bert	Seaquest	S. Invaders
Random	354	1.2	0	-20.4	157	110	179
Sarsa [3]	996	5.2	129	-19	614	665	271
Contingency [4]	1743	6	159	-17	960	723	268
DQN	4092	168	470	20	1952	1705	581
Human	7456	31	368	-3	18900	28010	3690

Figure: Average total reward for a fixed number of steps.

- o DQN source code: <https://github.com/deepmind/dqn>

DQN extensions I

- o Double DQN (DeepMind, 2016) [29]: Use separate networks to select best action and evaluate best action to reduce overestimation bias

$$\min_w L(w) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}} \left[\left(r + \gamma Q(s', \arg \max_{a'} Q(s', a'; w); w^-) - Q(s, a; w) \right)^2 \right]$$

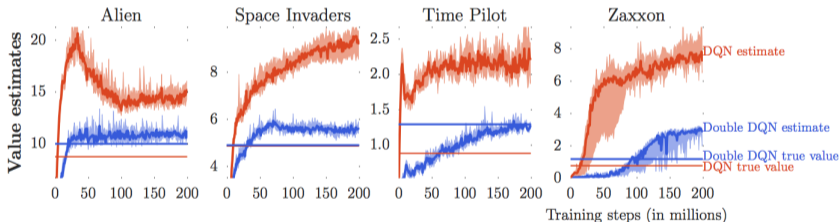
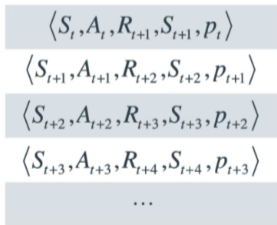


Figure: Value estimates by DQN (orange) and Double DQN (blue) on Atari games. The straight horizontal lines are computed by running the corresponding agents after learning concluded, and averaging the actual discounted return obtained from each visited state.

DQN extensions II

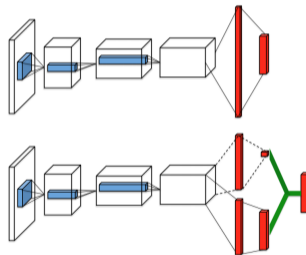
- **DQN with prioritized experience replay [21]**: Prioritize transitions in proportion to the absolute Bellman error

$$p \propto \left| r + \gamma \max_{a'} Q(s', a'; w) - Q(s, a; w) \right|$$



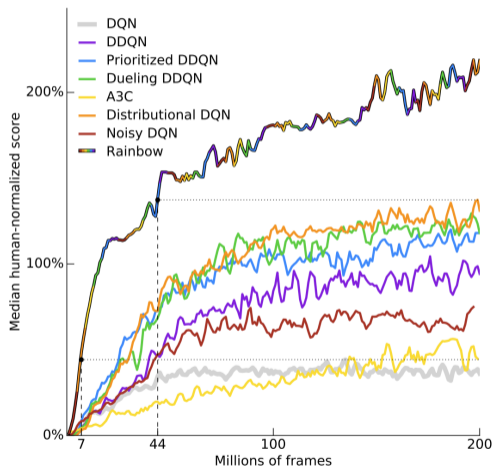
- **Dueling DQN [31]**: Split Q-networks into two streams to estimate value function and advantage function

$$Q(s, a; w, \alpha, \beta) = V(s; w, \beta) + \bar{A}(s, a; w, \alpha)$$



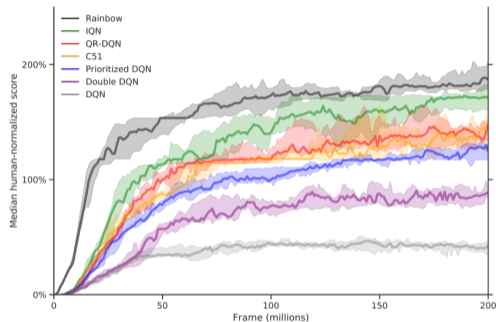
DQN mega extension

- Can these extensions be combined? Yes, Rainbow [8]!



The big zoo of DQN

Directory	Paper
dqn	Human Level Control Through Deep Reinforcement Learning
double_q	Deep Reinforcement Learning with Double Q-learning
prioritized	Prioritized Experience Replay
c51	A Distributional Perspective on Reinforcement Learning
qrdqn	Distributional Reinforcement Learning with Quantile Regression
rainbow	Rainbow: Combining Improvements in Deep Reinforcement Learning
iqn	Implicit Quantile Networks for Distributional Reinforcement Learning



Plot of median human-normalized score over all 57 Atari games for each agent

o Source code: https://github.com/deepmind/dqn_zoo

Policy-based/Actor-Critic DRL

- Combine the actor-critic approach with Deep Q Network
 - Asynchronous Advantage Actor-Critic (A3C) [13]
 - Soft Actor Critic (SAC) [7]
 - Deep deterministic policy gradient (DDPG) [12]: continuous control
 - Twin Delayed DDPG (TD3) [5]: continuous control
 -

A3C [13]

- o Idea: advantage actor-critic + deep Q-network + asynchronous implementation

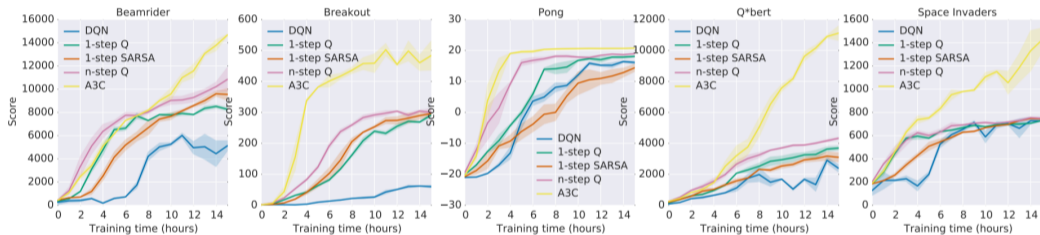


Figure: Comparison for DQN and A3C on five Atari 2600 games. 1-step Q means asynchronous one-step Q-learning.

DDPG [12] and TD3 [5]

- o **DDPG**: deterministic policy gradient + deep Q-network
 - o Select action $a \sim \mu(s; \theta) + \mathcal{N}(0, \sigma^2)$ (add noise to enhance exploration)
 - o Policy update: $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_i \nabla_a Q_w(s_i, \mu(s_i; \theta)) \nabla_{\theta} \mu(s_i; \theta)$
- o **TD3**: DDPG + clipped action exploration + delayed policy update + pessimistic double Q-learning
 - o Select action $a \sim \mu(s; \theta) + \epsilon$, $\epsilon \sim \text{clip}(N(0, \sigma^2), -c, c)$
 - o Delayed policy update: update critic more frequent than policy

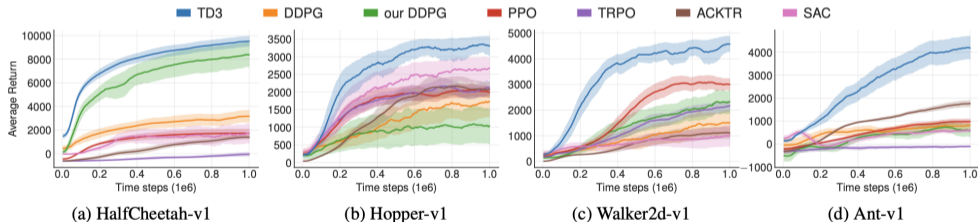


Figure: Learning curves for the OpenAI gym continuous control tasks.

Summary

o Deep Value-based Methods

- o DQN
- o Double DQN
- o Dueling DQN
- o DQN with prioritized experience replay
- o Rainbow
- o

o Deep Policy-based/Actor-Critic Methods

- o TRPO
- o PPO
- o A3C
- o SAC
- o DDPG/TD3
- o

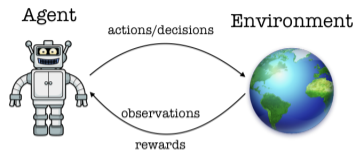
Question:

So, which one should we choose in practice? when do they work well?

Deep RL resources

- OpenAI Spinning up: <https://spinningup.openai.com/>
- The awesome list of deep RL (libraries and tutorials): <https://github.com/kengz/awesome-deep-rl>

Reinforcement learning



- Environment: Markov Decision Process (MDP) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, T, \gamma, \mu, r)$
- Agent: Parameterized deterministic policy $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$, where $\theta \in \Theta$

Reinforcement learning (RL) game

At time step $t = 0$: $S_0 \sim \mu(\cdot)$

for $t = 1, 2, \dots$ do:

agent observes the environment's state $S_t \in \mathcal{S}$

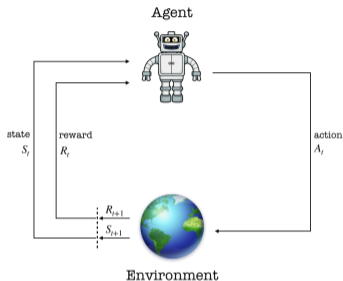
agent chooses an action $A_t = \pi_\theta(S_t) \in \mathcal{A}$

agent receives a reward $R_{t+1} = r(S_t, A_t)$

agent finds itself in a new state $S_{t+1} \sim T(\cdot | S_t, A_t)$

Exploration vs. exploitation in RL

- Challenge: Exploration vs. exploitation!



- Objective (non-concave):

$$\max_{\theta \in \Theta} J(\theta) := \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} R_t \mid \pi_{\theta}, \mathcal{M} \right]$$

- The environment only reveals the rewards after actions
- Exploitation: Maximize objective by choosing the appropriate action
- Exploration: Gather information on other actions

An optimization interpretation

- Objective (non-concave):

$$\max_{\theta \in \Theta} J(\theta) := \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} R_t \mid \pi_{\theta}, \mathcal{M} \right]$$

- Exploitation: Progress in the gradient direction

$$\theta_{t+1} \leftarrow \theta_t + \eta_t \nabla_{\theta} \widehat{J}(\theta_t)$$

- Exploration: Add stochasticity while collecting the episodes

- noise injection in the action space

[24, 12]

$$a = \pi_{\theta}(s) + \mathcal{N}(0, \sigma^2 I)$$

- noise injection in the parameter space

[19]

$$\tilde{\theta} = \theta + \mathcal{N}(0, \sigma^2 I)$$

Reinforcement learning with Langevin dynamics I

- Explore via an infinite dimensional concave-problem (linear in p):

$$\underset{p \in \mathcal{M}(\Theta)}{\text{maximize}} \quad \mathbb{E}_{\theta \sim p} [J(\theta)]$$

- $\mathcal{M}(\Theta)$ is the (infinite dimensional) space of all probability distributions on Θ .
- $p^* = \arg \max_p \mathbb{E}_{\theta \sim p} [J(\theta)]$ is a delta measure centered at $\theta^* = \arg \max_{\theta} J(\theta)$.

Reinforcement learning with Langevin dynamics II

- Exploit via a well-known entropy smoothing trick:

$$\underset{p \in \mathcal{M}(\Theta)}{\text{maximize}} \quad \mathbb{E}_{\theta \sim p} [J(\theta)] + \beta H(p)$$

- $H(p) = \mathbb{E}_{\theta \sim p} [-\log p(\theta)]$ is the entropy of the distribution p .
- the optimal solution takes the form $p_{\beta}^*(\theta) \propto \exp\left(\frac{1}{\beta} J(\theta)\right)$.
- Our proposal for explore-exploit
 - Use Langevin dynamics [32] to draw samples from $p_{\beta}^*(\theta)$
 - Use homotopy on the smoothing parameter β

Learning robust policies

- Why robust RL? In short: Generalization under environmental changes
 - upshots: self-driving car in varying environmental conditions
 - trends: from simple parametric models to super expressive neural networks
 - challenges: computational costs as well as the difficulty of training
- Highlight: Robust Adversarial Reinforcement Learning (RARL) [18]
 - train an **agent** neural net
 - train an **adversary** neural net
 - setup a minimax game between the two
- Several variants exist [16, 33]
 - Action Robust RL [28]

Two-Player Zero-Sum Markov Game

- Players:
 - Environment: Markov Decision Process (MDP) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \bar{\mathcal{A}}, T, \gamma, r, \mu)$
 - Agent: parameterized deterministic policy $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$, where $\theta \in \Theta$
 - Adversary: parameterized deterministic policy $\nu_\omega : \mathcal{S} \rightarrow \bar{\mathcal{A}}$, where $\omega \in \Omega$

Two-Player Zero-Sum Markov Game

At time step $t = 0$: $S_0 \sim \mu(\cdot)$

for $t = 1, 2, \dots$ do:

both players observe the environment's state $S_t \in \mathcal{S}$

both players choose the actions $A_t = \pi_\theta(S_t) \in \mathcal{A}$, and $\bar{A}_t = \nu_\omega(S_t) \in \bar{\mathcal{A}}$

the agent gets a reward $R_{t+1} = r(S_t, A_t, \bar{A}_t)$ while the adversary gets $-R_{t+1}$

both players find themselves in a new state $S_{t+1} \sim T(\cdot | S_t, A_t, \bar{A}_t)$

- Performance objective:

$$\max_{\theta \in \Theta} \min_{\omega \in \Omega} J(\theta, \omega) := \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} R_t \mid \pi_\theta, \nu_\omega, \mathcal{M} \right]$$

Robust Adversarial Reinforcement Learning (RARL)

- A natural *pure* strategy-based minimax objective

$$\max_{\theta \in \Theta} \min_{\omega \in \Omega} J(\theta, \omega).$$

- θ : an **agent** neural net
- ω : an **adversary** neural net
- *highly* non-concave/non-convex objective

- Theoretical challenges
 - a saddle point might NOT exist
 - no provably convergent algorithm

- Practical challenges
 - the simple (alternating) SGD does NOT work well in practice
 - adaptive methods (Adam, RMSProp,...) highly unstable, heavy tuning

[4]

RARL: From pure to mixed Nash Equilibrium

- Objective of RARL is a pure strategy formulation:

$$\max_{\theta \in \Theta} \min_{\omega \in \Omega} J(\theta, \omega).$$

- A new objective of RARL: Our **mixed** strategy proposal via game theory

$$\max_{p \in \mathcal{M}(\Theta)} \min_{q \in \mathcal{M}(\Omega)} \mathbf{E}_{\theta \sim p} \mathbf{E}_{\omega \sim q} [J(\theta, \omega)].$$

- where $\mathcal{M}(\mathcal{Z}) := \{\text{all (regular) probability measures on } \mathcal{Z}\}$.
- Existence of NE (p^*, q^*) : Glicksberg's existence theorem

[6].

A re-thinking of RARL via the mixed Nash equilibrium

- **Upshot:** Our mixed Nash Equilibrium proposal \equiv bi-linear matrix games

$$\begin{aligned} & \max_{p \in \mathcal{M}(\Theta)} \min_{q \in \mathcal{M}(\Omega)} \mathbf{E}_{\theta \sim p} \mathbf{E}_{\omega \sim q} [J(\theta, \omega)] \\ & \quad \Updownarrow \\ & \max_{p \in \mathcal{M}(\Theta)} \min_{q \in \mathcal{M}(\Omega)} \langle p, Gq \rangle \end{aligned}$$

- Caveat: **Infinite dimensions!!!**
- Key ingredients moving forward
 - $\langle p, h \rangle := \int h dp$ for a measure p and function h (Riesz representation)
 - the linear operator G and its adjoint G^\dagger :

$$\begin{aligned} (Gq)(\theta) &:= \mathbf{E}_{\omega \sim q} [J(\theta, \omega)] \\ (G^\dagger p)(\omega) &:= \mathbf{E}_{\theta \sim p} [J(\theta, \omega)], \end{aligned}$$

where $G : \mathcal{M}(\Omega) \rightarrow \mathcal{F}(\Theta)$, and $G^\dagger : \mathcal{M}(\Theta) \rightarrow \mathcal{F}(\Omega)$.

Training Phase

- We use the following special adversary with $\alpha = 0.1$ (Noisy Action Robust MDP):

Noisy Action Robust MDP Game

for $t = 1, 2, \dots$ do:

both players observe the environment's state $S_t \in \mathcal{S}$

both players choose the actions $A_t = \mu(S_t) \in \mathcal{A}$, and $A'_t = \nu(S_t) \in \mathcal{A}$

the resulting action $\bar{A}_t = (1 - \alpha)A_t + \alpha A'_t$ is executed in the environment \mathcal{M}

the agent gets a reward $R_{t+1} = r(S_t, \bar{A}_t)$ while the adversary gets $-R_{t+1}$

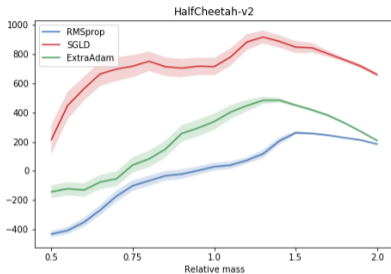
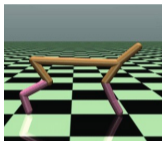
both players find themselves in a new state S_{t+1}

- We train the policy based on specific environment parameters
 - i.e., standard relative mass variables in OpenAI gym.

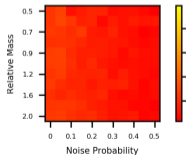
Testing Phase

- Robustness under Adversarial Disturbances (x-axis of the heatmap):
 - measure performance in the presence of an adversarial disturbance.
- Robustness to Test Conditions (y-axis of the heatmap):
 - measure performance with respect to varying test conditions.

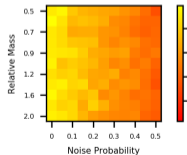
Experimental evaluation via MuJoCo



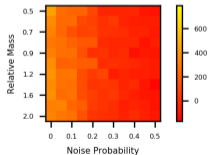
RMSprop/HalfCheetah-v2



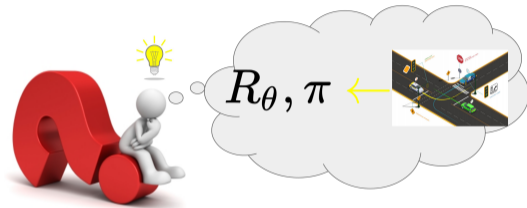
SGLD/HalfCheetah-v2



ExtraAdam/HalfCheetah-v2



A motivation for inverse reinforcement learning (IRL)



- The reward function is difficult to design in real world problems
- It is easier/more natural to use “demonstrations” by experts

The RL and IRL dichotomy

	IRL	RL
Input	Expert Demonstrations	Reward Function
Output	Optimal policy Reward function	Optimal Policy

- RL recovers a nearly optimal behavior from reward functions
- IRL recovers a nearly optimal behavior from demonstrations by an expert

Motivation for Robust IRL (This work)

- Mismatches between the settings of the expert and the learner
- Example: transfer the driving skills among different road conditions, traffic dynamics and car brands



Figure: A Toyota Prius¹ and Bugatti la voiture noir² have arguably different dynamics!

¹<https://www.autobild.de/artikel/toyota-prius-3-hybridauto-als-gebrauchtwagen-16425701.html>,

²<https://www.autobild.de/artikel/toyota-prius-3-hybridauto-als-gebrauchtwagen-16425701.html>

Basics: Markov Decision Processes (MDPs)

- A Markov Decision Process (MDP) is a tuple $(\mathcal{S}, \mathcal{A}, \gamma, T, r, \mu)$
 - \mathcal{S} is the state space
 - \mathcal{A} is the action space
 - $T : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$ is a mapping from state action pairs to distribution over the state space \mathcal{S}
 - γ is a scalar between 0 and 1 that is known as *discount factor*
 - $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is a mapping from state action pairs to a scalar value called *reward*
 - $\mu \in \Delta_{\mathcal{S}}$ is a probability distribution over states
- In particular, $T(s'|s, a)$ denotes the probability of landing in state s' after taking action a from state s

From policies to trajectories with a bit more notation

- A policy $\pi : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$ is a mapping from a state to a probability distribution over actions
- In the sequel, by a *trajectory*, we mean

$$\tau = (s_0, a_0, s_1, a_1, s_2, a_2, \dots).$$

- The probability of a trajectory factorizes as follows:

$$p_{\pi, T}(\tau) = \prod_{i=0}^{\infty} T(s_{i+1} | s_i, a_i) \pi(a_i | s_i) \mu(s_0).$$

Optimal policy vs optimal occupancy measure

- With the MDP formalism, RL solves the following problem

$$\max_{\pi \in \Delta_{\mathcal{A}}} \mathbf{E}_{\tau \sim p_{\pi, T}} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right] := V^{\pi}(s) \quad \forall s \in \mathcal{S}$$

- For IRL, we write the same objective as function of the *occupancy measure* λ :

$$\begin{aligned} \max_{\lambda} \langle \lambda, r \rangle &:= \sum_{s, a} \lambda(s, a) r(s, a) \\ \text{subject to} \quad &\sum_a \lambda(s, a) = \gamma \sum_{s', a'} T(s|s', a') \lambda(s', a') + (1 - \gamma) \mu(s) \quad \forall s \in \mathcal{S} \end{aligned}$$

- where the occupancy measure is the **discounted expected number of visits** for s, a

$$\lambda_T^{\pi}(s, a) = (1 - \gamma) \mathbb{E}_{p_{\pi, T}} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{1}_{((s_t, a_t) = (s, a))} \right].$$

Towards an IRL formulation

- Recall that the reward function r is unknown
- How can we learn from the demonstrations?
 - we can estimate the expert's occupancy measure λ_T^E from expert's trajectories
- Key Fact: $\forall \pi : \lambda_T^\pi = \lambda_T^E \implies \forall r \quad \langle \lambda_T^\pi, r \rangle = \langle \lambda_T^E, r \rangle$
 - A policy π matching the expert's occupancy measure results in the same performance!

IRL is a feasibility problem

- A simple feasibility problem

$$\begin{aligned} & \max_{\pi} 0 \\ & s.t. \quad \lambda_T^{\pi}(s, a) = \lambda_T^E(s, a) \quad \forall s, a \in \mathcal{S} \times \mathcal{A} \end{aligned}$$

- akin to moment-matching
- Pick one that maximizes the expected entropy of the policy π

$$\begin{aligned} & \max_{\pi} \sum_s \lambda_T^{\pi}(s) H^{\pi}(s) \\ & s.t. \quad \lambda_T^{\pi}(s, a) = \lambda_T^E(s, a) \quad \forall s, a \in \mathcal{S} \times \mathcal{A} \end{aligned}$$

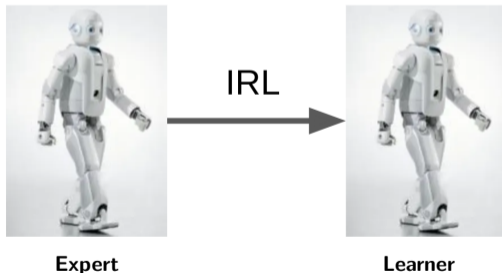
- $H^{\pi}(s) = - \sum_a \pi(a|s) \log \pi(a|s)$
- aka maximum causal entropy (MCE) IRL (already introduced in [36])
- also has a “dual” purpose

A critical limitation of MCE-IRL formulation

- Need the same dynamics T between expert and learner!

$$\max_{\pi} \sum_s \lambda_T^{\pi}(s) H^{\pi}(s)$$
$$s.t. \quad \lambda_T^{\pi}(s, a) = \lambda_T^E(s, a) \quad \forall s, a \in \mathcal{S} \times \mathcal{A}$$

- Factory-produced expert and learner setting is not realistic

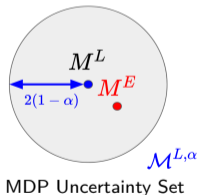


<https://www.autoevolution.com/news/samsung-steps-things-up-with-robtoray-walking-robot-video-50602.html>

Towards overcoming the limitation of MCE-IRL: Robust MCE IRL

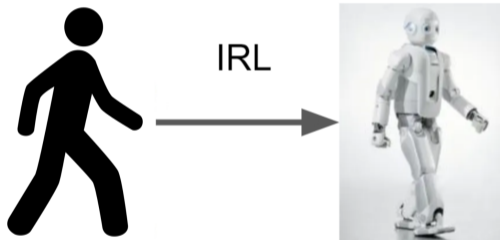
- In our work, we consider different transition dynamics:
 - T^E for the expert
 - T^L for the learner
- As a running example, We assume that T^L is within an uncertainty set centered around: T^L , i.e.,

$$\mathcal{T}^{L,\alpha} = \{T = \alpha T^L + (1 - \alpha)\bar{T} \quad \forall \bar{T} \in \Delta_S\}$$



On the generality of the uncertainty set

- For $\alpha = 0$, the set $\mathcal{T}^{L,\alpha}$ can correspond to **any possible transition dynamics** \bar{T}
- Example: The expert can be a human



Expert

Learner

https://commons.wikimedia.org/wiki/File:Man_walking_icon_1410105361.svg

Robust MCE IRL formulation

- Added twist: the MCE IRL with an additional minimization over the uncertainty set

$$\begin{aligned} \max_{\pi} \min_{T \in \mathcal{T}^{L, \alpha}} \sum_s \lambda_T^{\pi}(s) H^{\pi}(s) \\ \text{s.t. } \lambda_T^{\pi}(s, a) = \lambda_{TE}^E(s, a) \quad \forall s, a \in \mathcal{S} \times \mathcal{A} \end{aligned}$$

- We leverage Lagrangian duality for the numerical solutions

Lagrangian of robust MCE IRL

- Introduce the Lagrangian

$$\min_r \max_{\pi} \min_{T \in \mathcal{T}^{L, \alpha}} \sum_s \lambda_T^\pi(s) H^\pi(s) + \langle r, \lambda_T^\pi - \lambda_{TE}^E \rangle \quad \text{with } r \in \mathbb{R}^{S \times \mathcal{A}}$$

- the dual variable acts as the unknown reward r
- Via Danskin's theorem, compute the gradients to update the dual with a step-size η at iteration k :

$$r_{k+1} \leftarrow r_k - \eta(\lambda_{T_k}^{\pi_k} - \lambda_{TE}^E) \quad (\text{Reward Update})$$

- Then, (π_k, T_k) is a saddle point of the following min-max problem

$$\max_{\pi} \min_{T \in \mathcal{T}^{L, \alpha}} \sum_s \lambda_T^\pi(s) H^\pi(s) + \langle r_k, \lambda_T^\pi \rangle = \max_{\pi} \min_{T \in \mathcal{T}^{L, \alpha}} \langle r_k + H^\pi(s), \lambda_T^\pi \rangle \quad (\text{Robust MDP})$$

- Can use more sophisticated methodology but this one is something we can analyze

Subtleties on the policy update step

- In Robust MDP, $r_k + H^\pi(s)$ acts as a reward function in the policy update step
- We leverage the idea of solving Robust MDPs via a zero sum Markov games [28, 10], i.e., solving:

$$\max_{\pi} \min_{\pi^{\text{op}}} \langle r_k + H^\pi(s), \lambda_{T^L}^{\alpha\pi + (1-\alpha)\pi^{\text{op}}} \rangle$$

- can be solved sampling trajectories only with T^L
- more efficient than solving the min over the MDP uncertainty set $\mathcal{T}^{L,\alpha}$
- the latter would require to sample trajectories from any environment in $\mathcal{T}^{L,\alpha}$

The algorithm

Algorithm 1 Robust MCE IRL via Markov Game

Input: opponent strength $1 - \alpha$

Initialize: player policy π_0^{Pl} , opponent policy π_0^{OP} , and initial reward parameters r .

while not converged **do**

- Compute $\rho_{ML}^{\alpha\pi_k^{\text{Pl}} + (1-\alpha)\pi_k^{\text{OP}}}$ by dynamic programming as in MCE IRL(see [2]).
- Update reward:

$$r_{k+1} \leftarrow r_k - \left(\lambda_{ML}^{\alpha\pi_k^{\text{Pl}} + (1-\alpha)\pi_k^{\text{OP}}} - \lambda_{TE}^E \right) \quad (\text{Reward Update})$$

- Fix the reward r_{k+1} to update π^{Pl} and π^{OP} s.t. they solve the problem.

$$(\pi_{k+1}^{\text{Pl}}, \pi_{k+1}^{\text{OP}}) = \max_{\pi} \min_{\pi^{\text{OP}}} \langle r_{k+1} + H^{\pi}(s), \lambda_{TL}^{\alpha\pi + (1-\alpha)\pi^{\text{OP}}} \rangle \quad (\text{Solve Zero Sum Game})$$

end while

Output: player policy π^{Pl}

Theoretical guarantees

Theorem (Stylized version of Theorem 9 [30])

For any parameter of the MDP uncertainty set α , let us assume

- $|r(s, a)| \leq R \quad \forall s, a \in \mathcal{S} \times \mathcal{A}$
- $R = (1 - \gamma)^2$
- The expert dynamics T_E minimizes Robust MDP
- $d_{\text{dyn}}(T^L, T^E) = \max_{s, a} \|T^L(\cdot|s, a) - T^E(\cdot|s, a)\|_1 \leq 1$

Then, we have the following bound for the performance in the learner environment T^L :

$$\max_{\pi} V_{T^L}^{\pi} - V_{T^L}^{\pi^{\text{Pl}}} \leq d_{\text{dyn}}(T^L, T^E) + 2 \left[(1 - \alpha)^2 + \alpha \cdot d_{\text{dyn}}(T^E, T^L) \right]$$

Comparison of MCE-IRL and Robust MCE-IRL

Corollary (To the stylized version of Theorem 9 [30])

It follows that, for $\alpha = 1$ (MCE-IRL), using the notation $V_{T^L}^{\text{MCE}} = V_{T^L}^{\pi^{\text{pl}}}$, we have

$$\max_{\pi} V_{T^L}^{\pi} - V_{T^L}^{\text{MCE}} \leq 3d_{\text{dyn}}(T^L, T^E).$$

For the optimal tuning of $\alpha = 1 - \frac{d_{\text{dyn}}(T^L, T^E)}{2}$, using the notation, we instead have $V_{T^L}^{\text{Robust}} = V_{T^L}^{\pi^{\text{pl}}}$

$$\max_{\pi} V_{T^L}^{\pi} - V_{T^L}^{\text{Robust}} \leq 3d_{\text{dyn}}(T^L, T^E) - \frac{(d_{\text{dyn}}(T^L, T^E))^2}{2}.$$

Remark: o The paper presents a constructive example for which the bound holds with equality for MCE-IRL.

A simple demonstration

- We test our algorithm in a Gridworld problem
- The agent starts from a state drawn uniformly at random state
- The goal is to reach the top left corner where the reward is non negative

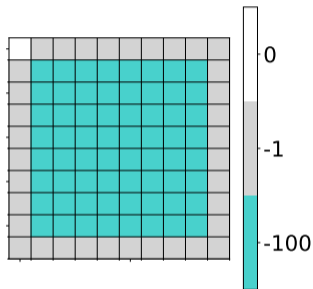


Figure: Schematics of the environment. Gridworld environment

- We introduce the variable *Expert Noise* as ϵ^E , and define the expert dynamics as follows:

$$T^E(s'|s, a) = (1 - \epsilon^E)T^L(s'|s, a) + \epsilon^E U(s'|s, a) \quad \forall s', s, a \in \mathcal{S} \times \mathcal{S} \times \mathcal{A}$$

where $U(\cdot|s, a)$ is a uniform distribution over the states that are first neighbors of s .

Effect of the noise

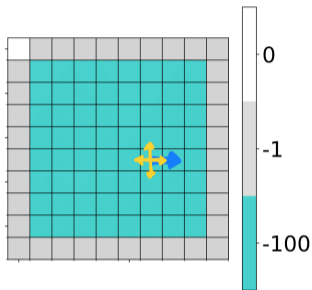


Figure: Effect of the noise on the Gridworld environment when the agent selects the action up (blue arrow).

- In this particular example, the agent takes action RIGHT and T^L is deterministic.
- With probability $1 - \epsilon^E$, the agent follows the blue arrow.
- With probability ϵ^E , it moves according to the yellow arrows.
- The noise is proportional to the mismatch, i.e., $\epsilon^E = \frac{d_{\text{dyn}}(T^L, T^E)}{2} \left(1 - \frac{1}{|\mathcal{S}|}\right)$

Results

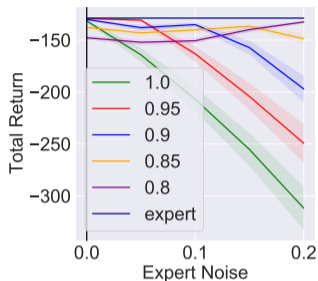


Figure: On the **x-axis** we report the noise in the expert environment ϵ^E . On the **y-axis** we have the performance in the learner environment. The **legend** contains the different values of α .

In the tabular experiments, we notice the following:

- V_{TL}^{MCE} , the green line, decays as the expert noise increases.
- The other lines represent V_{TL}^{Robust} for different values of α .
- In agreement with the theory the choice $\alpha = 1 - \epsilon^E$ performs the best.

Function approximation for continuous state and action pairs

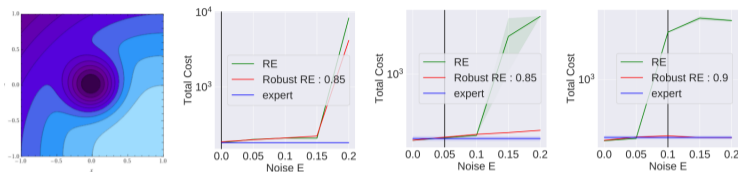


Figure: Linear function approximation. On the **x-axis** we report the noise in the expert environment. On the **y-axis** we have the performance in the learner environment. The **legend** contains the different values of α . The black vertical line denotes the noise in the learner environment.

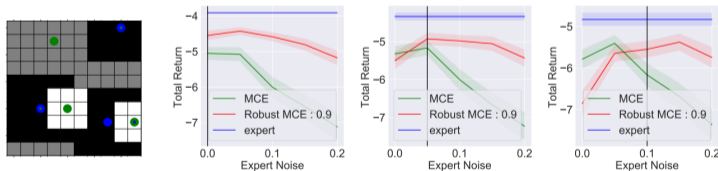


Figure: Nonlinear function approximation. On the **x-axis** we report the noise in the expert environment. On the **y-axis** we have the performance in the learner environment. The **legend** contains the different values of α . The black vertical line denotes the noise in the learner environment.

Conclusions

- Robust formulation of MCE IRL & an efficient solution
- Encouraging theoretical analysis showing provable improvements if α is chosen appropriately
- Numerical evidence corroborating the performance claims

References I

- [1] Amir Beck and Marc Teboulle.
Mirror descent and nonlinear projected subgradient methods for convex optimization.
Operations Research Letters, 31(3):167–175, 2003.
74
- [2] Michael Bloem and Nicholas Bambos.
Infinite time horizon maximum causal entropy inverse reinforcement learning.
In *53rd IEEE Conference on Decision and Control*, pages 4911–4916, 2014.
58
- [3] Vivek S Borkar and Vijaymohan R Konda.
The actor-critic algorithm as multi-time-scale stochastic approximation.
Sadhana, 22(4):525–543, 1997.
11
- [4] Partha Dasgupta and Eric Maskin.
The existence of equilibrium in discontinuous economic games, i: Theory.
The Review of economic studies, 53(1):1–26, 1986.
38
- [5] Scott Fujimoto, Herke Hoof, and David Meger.
Addressing function approximation error in actor-critic methods.
In *International Conference on Machine Learning*, pages 1582–1591, 2018.
26, 28

References II

- [6] Irving L Glicksberg.
A further generalization of the kakutani fixed point theorem, with application to nash equilibrium points.
Proceedings of the American Mathematical Society, 3(1):170–174, 1952.
39
- [7] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine.
Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor.
In *International Conference on Machine Learning*, pages 1856–1865, 2018.
10, 26
- [8] Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver.
Rainbow: Combining improvements in deep reinforcement learning.
In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
24
- [9] S. Kakade.
A natural policy gradient.
In *Advances in Neural Information Processing Systems (NeurIPS)*, 2001.
10
- [10] Parameswaran Kamalaruban, Yu-Ting Huang, Ya-Ping Hsieh, Paul Rolland, Cheng Shi, and Volkan Cevher.
Robust reinforcement learning via adversarial training with langevin dynamics.
In *NeurIPS*, 2020.
57

References III

- [11] Vijay Konda and John Tsitsiklis.
Actor-critic algorithms.
Advances in neural information processing systems, 12, 1999.
11
- [12] 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.
arXiv preprint arXiv:1509.02971, 2015.
26, 28, 33
- [13] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu.
Asynchronous methods for deep reinforcement learning.
In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.
26, 27
- [14] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.
Human-level control through deep reinforcement learning.
Nature, 518(7540):529–533, 2015.
16, 19, 21

References IV

[15] Arkadi Nemirovski.

Prox-method with rate of convergence $\mathcal{O}(1/t)$ for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems.

SIAM Journal on Optimization, 15(1):229–251, 2004.

74

[16] Arnab Nilim and Laurent El Ghaoui.

Robust control of Markov decision processes with uncertain transition matrices.

Operations Research, 53(5):780–798, 2005.

36

[17] Jan Peters and Stefan Schaal.

Natural actor-critic.

Neurocomputing, 71(7-9):1180–1190, 2008.

10

[18] Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta.

Robust adversarial reinforcement learning.

In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2817–2826. JMLR.org, 2017.

36

[19] Matthias Plappert, Rein Houthoofd, Prafulla Dhariwal, Szymon Sidor, Richard Y Chen, Xi Chen, Tamim Asfour, Pieter Abbeel, and Marcin Andrychowicz.

Parameter space noise for exploration.

arXiv preprint arXiv:1706.01905, 2017.

33

References V

- [20] Itay Safran, Ronen Eldan, and Ohad Shamir.
Depth separations in neural networks: What is actually being separated?
In Alina Beygelzimer and Daniel Hsu, editors, *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 2664–2666. PMLR, 25–28 Jun 2019.
15
- [21] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver.
Prioritized experience replay.
arXiv preprint arXiv:1511.05952, 2015.
23
- [22] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz.
Trust region policy optimization.
In *International conference on machine learning*, pages 1889–1897. PMLR, 2015.
10
- [23] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel.
High-dimensional continuous control using generalized advantage estimation.
arXiv preprint arXiv:1506.02438, 2015.
10
- [24] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller.
Deterministic policy gradient algorithms.
In *ICML*, 2014.
33

References VI

- [25] Richard S Sutton and Andrew G Barto.
Reinforcement learning: An introduction.
MIT press, 2018.
4
- [26] Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour.
Policy gradient methods for reinforcement learning with function approximation.
In *Advances in neural information processing systems*, pages 1057–1063, 2000.
7
- [27] Matus Telgarsky.
Benefits of depth in neural networks.
In *Conference on learning theory*, pages 1517–1539. PMLR, 2016.
15
- [28] Chen Tessler, Yonathan Efroni, and Shie Mannor.
Action robust reinforcement learning and applications in continuous control.
In *International Conference on Machine Learning*, pages 6215–6224. PMLR, 2019.
36, 57
- [29] Hado Van Hasselt, Arthur Guez, and David Silver.
Deep reinforcement learning with double q-learning.
In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
22

References VII

- [30] Luca Viano, Yu-Ting Huang, Parameswaran Kamalaruban, Adrian Weller, and Volkan Cevher.
Robust inverse reinforcement learning under transition dynamics mismatch, 2021.
59, 60
- [31] Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas.
Dueling network architectures for deep reinforcement learning.
In *International Conference on Machine Learning*, pages 1995–2003, 2016.
23
- [32] Max Welling and Yee W Teh.
Bayesian learning via stochastic gradient langevin dynamics.
In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 681–688, 2011.
35
- [33] Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem.
Robust markov decision processes.
Mathematics of Operations Research, 38(1):153–183, 2013.
36
- [34] Yue Wu, Weitong Zhang, Pan Xu, and Quanquan Gu.
A finite time analysis of two time-scale actor critic methods.
arXiv preprint arXiv:2005.01350, 2020.
11

References VIII

- [35] Dmitry Yarotsky.
Error bounds for approximations with deep relu networks.
Neural Networks, 94:103–114, 2017.
15
- [36] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al.
Maximum entropy inverse reinforcement learning.
volume 8, pages 1433–1438. Chicago, IL, USA, 2008.
51

Supplementary: Entropic mirror descent iterates in infinite dimension

- Negative Shannon entropy and its Fenchel dual: ($dz := \text{Lebesgue}$)
 - $\Phi(p) = \int p \log \frac{dp}{dz}$.
 - $\Phi^*(h) = \log \int e^h$.
 - $d\Phi$ and $d\Phi^*$: Fréchet derivatives.³

Theorem (Infinite-dimensional mirror descent, informal)

For a learning rate η , a probability measure p , and an arbitrary function h , we can equivalently define

$$p_+ = \mathbf{MD}(p, h) \quad \equiv \quad p_+ = d\Phi^*(d\Phi(p) - \eta h) \quad \equiv \quad dp_+ = \frac{e^{-\eta h} dp}{\int e^{-\eta h} dp}.$$

Moreover, most the essential ingredients in the analysis of finite-dimensional prox methods can be generalized to infinite dimension.

- Continuous analog of the entropic mirror descent
- Mirror-prox also possible

[1]

[15]

³Under mild regularity conditions on the measure/function.

Supplementary: Entropic mirror descent in infinite dimension: rates

o Algorithm:

Algorithm 1 Infinite-Dimensional Entropic Mirror Descent

Input: Initial distributions p_1, q_1 , and learning rate η

for $t = 1, 2, \dots, T - 1$ **do**

$$p_{t+1} = \text{MD}_\eta(p_t, -Gq_t)$$

$$q_{t+1} = \text{MD}_\eta(p_t, G^\dagger p_t)$$

end for

Output: $\bar{p}_T = \frac{1}{T} \sum_{t=1}^T p_t$ and $\bar{q}_T = \frac{1}{T} \sum_{t=1}^T q_t$

Theorem (Convergence Rates)

Let $\Phi(p) = \int dp \log \frac{dp}{dz}$. Then

1. Entropic MD $\Rightarrow O(T^{-\frac{1}{2}})$ -NE.
2. If only stochastic derivatives ($\hat{G}^\dagger p$ and $-\hat{G}q$) are available, then Entropic MD $\Rightarrow O(T^{-\frac{1}{2}})$ -NE in expectation.