

Challenges of Real-World RL

Gabriel Dulac-Arnold¹, Daniel Mankowitz^{*}, Todd Hester^{*}

¹Google Research, Brain Team ^{*}Deepmind

Summary

- Reinforcement learning has proven its worth in numerous artificial domains.
- Many of the recent advances in RL research are hard to leverage in real-world systems, primarily due to rarely satisfied assumptions.
- We propose 9 challenges that must be addressed to bring RL to real-world systems, and discuss current approaches to each of these.
- We propose example domains, as well as a full training/evaluation regime that more closely resembles real-world systems and acts as a more realistic testbed for candidate RL algorithms.

Challenges

Real-world policies need to be:

- Off-policy & off-line:** policies need to be trained off-line from the fixed logs of an external behavior policy.
- Efficient:** Learning on the real system from limited samples.
- Scalable:** policies need to reason with high-dimensional continuous state and action spaces.
- Safe:** environments have safety constraints that should never or at least rarely be violated.
- Risk-Adverse & Robust:** environments may be partially observable, non-stationary or stochastic and potentially adversarial.
- Discerning:** policies need to reason with reward functions that are unspecified, multi-objective, or risk-sensitive.
- Explainable:** system operators require explainable policies and actions.
- Fast:** Inference must happen in real-time at the control frequency of the system.
- Mnemonic:** Policies should be able to handle large and/or unknown delays in the system actuators, sensors, or rewards.

Environments

Cart-Pole Variables: x, θ

Type	Constraint
Static	Limit range: $x_l < x < x_r$
Kinematic	Limit velocity near goal: $ \theta_c - \theta > \theta_L \vee \dot{\theta} < \dot{\theta}_V$
Dynamic	Limit cart acceleration: $\ddot{x} < A_{\max}$

Walker Variables: θ, u, F

Type	Constraint
Static	Limit joint angles: $\theta_L < \theta < \theta_U$ Enforce uprightness: $0 < u_x$
Kinematic	Limit joint velocities: $\max_i \dot{\theta}_i < L_{\dot{\theta}}$
Dynamic	Limit foot contact forces: $F_{\text{foot}} < F_{\max}$

Robustness

Env.	Noise	Non-Stationarity
Cart-Pole	Actuator and sensor delays	Track friction increasing with time
Walker	Noisy perception of terrain	Occasionally non-responsive leg actuator
Manipulator	Imprecise proprioception	Changes in gripper friction
Humanoid	Reduced torque on leg actuator	Varying payload CoGs

Manipulator Variables: θ, F, \mathcal{M}

Type	Constraint
Static	Limit joint angles $\theta_L < \theta < \theta_U$ Avoid dynamic obstacles $\mathcal{M} \cap \mathcal{M}_{O,i} = \emptyset$ Avoid self-contact $\mathcal{M} \cap \mathcal{M} = \mathcal{M}$
Kinematic	Limit joint velocities: $\max_i \dot{\theta}_i < L_{\dot{\theta}}$
Dynamic	Acceleration Limits: $\max_i \ddot{\theta}_i < L_{\ddot{\theta}}$ Limit end effector forces: $F_{\text{EE}} < F_{\max}$

Humanoid Variables: θ, u, F

Type	Constraint
Static	Limit joint angles: $\theta_{L,i} < \theta_i < \theta_{U,i}$ Enforce uprightness: $0 < u_x$
Kinematic	Limit joint velocities: $\max_i \dot{\theta}_i < L_{\dot{\theta}}$
Dynamic	Limit foot contact forces: $F_{\text{foot}} < F_{\max}$ Encourage falls on posterior $F_i < F_{\max,1} \forall i \in \mathcal{C} \setminus i_{\text{post}}$ $F_{\text{post}} < F_{\max,2}$

Evaluators

Challenge	Evaluator
Off-line	$J^{\text{start}} = R(\text{Train}(D_{\pi_B}))$
Efficient	$J^{\text{eff.}} = \min \mathcal{D}_i \text{ s.t. } R(\text{Train}(D_i)) > R_{\min}$
Safe	$J^{\text{safety}}(\pi) = \left(\sum_{i=1}^T c_j(s_i, a_i) \right)_{1 \leq j \leq K} \in \mathbb{R}^K$
Robust	$J^{\text{robust}}(\pi) = \frac{1}{K} \sum_{p \in \mathbf{P}} \mathbf{E}^p \left[\sum_{i=1}^T r(s_i, a_i) \right]$
Discerning	$J^{\text{multi}}(\pi) = \left(\sum_{i=1}^{T_n} r_j(s_i, a_i) \right)_{1 \leq j \leq K} \in \mathbb{R}^K$

Proposed Framework

- Training performed in Batch-RL (off-line, Initially off-policy)
 - Training from behavior policy logs with varying sizes and policy qualities.
- Training considered as on-line:
 - Use both safety conscious and robustness-enhanced environments.
 - Every interaction should be considered as the real system being run (goes towards efficiency counts and safety constraints).
- Evaluate according to the proposed evaluators, and compare training algorithms according to a multi-dimensional approach.

