



Optimizing Decision Strategies through Advanced Learning Techniques

*Pradeep Kumar Reddy Dasari Leela

Corresponding author Email: pdasari1@umbc.edu

Abstract: This study explores an approach to refining decision- making models using adaptive learning methods. By structure data-driven strategies, the research enhances efficiency in complex scenarios. Various techniques are examined to improve adaptability and performance, with an emphasis on practical applications. The findings highlight key advancements that con- tribute to more effective and controlled learning processes.

1. INTRODUCTON

This section introduces fundamental concepts in Reinforce- ment Learning (RL), a framework where agents learn to make decisions by interacting with an environment.

A. Core Concepts

Consider a discrete-time system, where at each time step t, the environment is in a state st and the agent selects an action at. The agent's behaviour is determined by a policy $\pi\theta$, which can be deterministic (at = $\pi\theta$ (st)) or stochastic ($\pi\theta$ (at|st)). When states are only partially observable, the policy may depend on observations ot. The environment evolves based on a transition function or model, defined as a conditional probability distribution p (st+1|st, at). A sequence of states and actions forms a trajectory $\tau = (s1, a1... sT, at)$. The agent receives feedback through a reward function r(s, a), which can be known or learned (as in Inverse RL). The goal is to maximize the expected cumulative reward:



 $\theta = \arg \max E \sim "\Sigma r (st, at) #,$



FIGURE 2. Three-step reinforcement learning pipeline.

Where the trajectory distribution is:

 $P(\tau) = p(s) Y \pi (a | s) p (s | s, a).$

2) Value Function

The value function estimates the expected future reward I often assume the Markov property, where st+1 depends T only on st and at (see Fig. 1).

 $V \pi (st) =$ $E\pi[r (set', at') | st], t'=t$

To facilitate policy optimization, two key functions are defined:

1) **Q** Function

The Q-function quantifies the expected future reward from taking action at in state st under policy π :

 $T Q\pi$ (st, at) = $E\pi[r(st', at') | st, at]$. t'=t

With the relation V π (st) = Ea $\sim \pi$ (•|s) [Q π (st, at)]. The over- all objective becomes Es $\sim p(s)$ [V π (s1)].

C. RL Workflow

RL involves three key stages: (1) data generation through agent-environment interaction, (2) function approximation us- ing collected samples, and (3) policy improvement based on estimated rewards. This loop continues iteratively to enhance agent performance (see Fig. 2).

2. IMITATION LEARNING

Imitation learning, also known as behavioural cloning, aims to train an agent to replicate expert behaviour from demonstrate. The agent is trained to fit a policy π (at|ot) that maps observations to actions.

Distribution Mismatch

A core challenge is the distribution mismatch between training and test time. Small mistakes can lead the agent into unfamiliar states, compounding errors and deviating from the expert trajectory, as illustrated in Fig. 3.



FIGURE 3. Compounding errors in behaviour cloning.

Dataset Aggregation (D Agger)

To mitigate error accumulation, Dataset Aggregation (DAg- ger) [1] iteratively augments the training dataset with agent- collected states, relabeled by a human expert.

Algorithm 1 Dataset Aggregation (DAgger)

- 1. Require: Expert dataset $D = \{(oi, ai)\}$ N
- 2. while not converged do
- 3. Train policy $\pi\theta$ on D
- 4. Run $\pi\theta$ to collect observations DP
- 5. Reliable DP with expert actions
- 6. Aggregate: $D \leftarrow D \cup DP$
- 7. return Final policy $\pi \theta$

D Agger aligns the training and test distributions but relies heavily on human supervision, which can introduce noise and is not always Markova.

Failure Modes of Imitation Learning

1) Non-Markova Behaviour

Human decisions often depend on historical context, leading to a mismatch when fitting Markovian policies. One approach is to use recurrent neural networks (e.g., LSTM) to encode history, as shown in Fig. 4. However, longer histories can introduce causal confusion, where the agent learns spurious correlations. For example, an agent might associate braking with a dashboard light rather than an obstacle.



Typically, LSTM cells work better here FIGURE 4. Addressing non-Markovian data with RNN.

2) Multimodal Expert Behavior

Experts may exhibit multiple valid behaviors in the same state (e.g., avoiding a tree by steering left or right). A unimodal policy (e.g., Gaussian) may average actions, leading to unsafe behavior (e.g., going straight into the tree). Solutions include using mixtures of Gaussians, latent variable models, or autoregressive discretization.



FIGURE 5. Averaging multimodal actions can result in failure.

Theoretical Error Analysis

Let π * be the expert policy and assume the imitation policy $\pi\theta$ has error rate ϵ on training distribution ptrain(s). The state distribution under $\pi\theta$ becomes:

 $\begin{array}{l} P\theta \left(st \right) = \left(1 - \varepsilon \right) \text{ tptrain} \left(st \right) + \left(1 - \left(1 - \varepsilon \right) \text{ t} \right) \text{ pmistake} \left(st \right) \\ \text{The total variation divergence yields:} \\ \left| p\theta \left(st \right) - p \text{train} \left(st \right) \right| \leq 2\varepsilon t \\ \text{Thus, the expected number of mistakes grows as:} \\ \text{E} \left[\# \text{ mistakes } 2 \right] \in O \left(\varepsilon T \right) \text{ Imitation learning faces challenges such as distribution} \\ \text{Mismatch, non-Markovian data, multimodal behaviour, and compounding errors. While DAgger addresses some of these issues, it depends on human supervision, which may be noisy or insufficient for certain state-action spaces. Future work explores learning from unlabelled or synthetic data to overcome these limitations. \end{array}$

3. POLICY GRADIENT METHODS

Reinforcement Learning (RL) can be framed as an optimization problem: $\theta = \arg \max E \tau \sim p \theta(\tau) "\Sigma r(st, at) #$ Define the objective function: $J(\theta) = E\tau \sim \pi \theta(\tau) [r(\tau)] = \int \pi \theta(\tau) r(\tau) d\tau$ **Policy Gradient Theorem** Using the identity $\nabla \theta \pi \theta$ (τ) = $\pi \theta$ (τ) $\nabla \theta \log \pi \theta$ (τ), I get: $\nabla \theta J(\theta) = E\tau \sim \pi \theta(\tau) \left[\nabla \theta \log \pi \theta(\tau) r(\tau) \right]$ For trajectory $\tau = (s, a, s, a)$: $\text{Log } \pi\theta (\tau) = \log p (s1) + \Sigma \log \pi\theta (at|st) + \Sigma \log p (st+1|st, at)$ Only $\pi\theta$ (at|st) depends on θ , so: $\Theta \tau \sim \pi \theta (\tau) \theta t = 1 \theta t t$ t=1**Monte Carlo Approximation** Since the expectation is intractable, approximate using Monte Carlo with N samples: $\nabla \theta J(\theta) \approx N \Sigma T \nabla \theta \log \pi \theta$ (ai, t|si, and t)! ΣTR (si,t, ai,t) Update policy via gradient ascent: $\Theta \leftarrow \theta + \alpha \nabla \theta J(\theta)$ **Example: Gaussian Policy** For $\pi\theta$ (at|st) = N (f θ (st), Σ): $\text{Log } \pi\theta (\text{at}|\text{st}) = -2f\theta (\text{st}) - \text{at} \quad 2 + C\nabla$ Gradient: $\nabla \log \pi$ (a |s) = $-\Sigma - 1$ (f (s) – a) df θ **Algorithm 2 REINFORCE Algorithm** Require: Policy $\pi\theta$ (a|s), learning rate α While not converged do Sample N trajectories $\{\tau i\}$ from $\pi \theta$



FIGURE 6. Higher-reward trajectories become more probable

Intuition behind Policy Gradient

Policy gradient increases the probability of high-reward trajectories: $\nabla J(\theta) \approx 1 \Sigma \nabla \log \pi(\tau) r(\tau)$

Compared to Maximum Likelihood:

Policy gradient assigns more weight to high-reward trajectories, effectively learning to prefer them.

Policy Gradient in POMDPs

The policy gradient does not rely on the Markov property.

In POMDPs, replace states stwith observations ot:

 $\nabla \theta J(\theta) \approx Ni=1 \nabla \theta \log \pi \theta (ai,t|oi,t) t r(oi,t, ai,t)$

Variance Reduction Using Baselines

Causality

By causality, actions at time t cannot influence past rewards:

 $\nabla \theta J(\theta) \approx N$ $\Sigma \nabla \theta \log \pi \theta(ai,t|si,t) \Sigma Tr(si,t',ai,t')! i=1 t=1$

This "reward-to-go" reduces variance by excluding irrele- vant past rewards. $\theta \ \theta \ t \ t \theta \ tt \ d\theta$

Drawback of Naive Policy Gradient

Naive policy gradient is sensitive to reward scaling. Adding a constant to all rewards shifts their values but affects the gradient direction, leading to high variance. Variance-reducing to stabilize learning.

4. BASELINES IN POLICY GRADIENT

To reduce variance in policy gradient methods, baselines are introduced—typically the average reward—so that only above- average trajectories contribute to learning. Let the baseline be defined as:

 $b = 1 \Sigma r (\tau)$

Incorporating this into the policy gradient yields:

 $\nabla \; J \; (\theta) \approx 1 \; \Sigma \; \nabla \; log \; \pi \; (\tau) \; [r \; (\tau) - b \;$

5. ACTOR-CRITIC ALGORITHMS

In the previous chapter, I derived the policy gradient theorem

$$\nabla g_{i}(\mathcal{O}) \simeq \frac{1}{N} \sum_{t=1}^{N} \nabla_{\theta} \log \pi_{\theta}(g_{it}|s_{it})$$

$$\times \sum_{t'=t}^{t} r(s_{it}, g_{it}) \qquad (1)$$

where the term Q^{i} , trepresents the reward-to-go, i.e., the expected cumulative reward from time t onward. While this Monte Carlo-based estimate is unbiased, it suffers from high variance. I now explore how actor-critic methods improve upon this by leveraging better approximations of Q^{i} ,t.

Reward-to-Go and Q-Function

A more precise form of the reward-to-go is the expected return conditioned on a state-action pair:

$$\Omega(\mathbf{s}_t, a_t) = \sum_{t'=t}^{t} \mathbb{E}\left[f(\mathbf{s}_t', a_{t'})|\mathbf{s}_t, a_t\right]$$

Replacing Monte Carlo returns with Q (st, at) yields a refined policy gradient estimate:

$$\nabla \mathfrak{Q}_{\mathcal{A}}(\vartheta) \simeq \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\vartheta} \log \pi \mathfrak{Q}(\mathfrak{g}_{i,t}|\mathfrak{g}_{i,t}) Q(\mathfrak{g}_{i,t},\mathfrak{g}_{i,t})$$

Variance Reduction with Baselines

To further reduce variance, I subtract a baseline from Q. A common choice is the value function:

 $V(s_t) = \mathbb{E}_{q_t \sim \pi_0}[Q(s_t, a_t)]$

This gives rise to the advantage function, which quantifies the relative value of taking action at in state st:

$$\underline{A}^{\pi}(\underline{s}_{t}, a_{t}) = Q^{\pi}(\underline{s}_{t}, a_{t}) - V^{\pi}(\underline{s}_{t})$$

Substituting this into the policy gradient, I obtain:

$$\nabla \vartheta J(\vartheta) \simeq \frac{1}{N} \sum_{i=1}^{N} \nabla \vartheta \log \pi \vartheta (a i \mathbf{s} | s i \mathbf{s}) A^{\pi}(s i \mathbf{s} | s i \mathbf{s})$$

Approximating the Advantage Function

To approximate $A\pi$, I leverage the relation $Q^{\pi}(s_t, a_t) = r(s_t, a_t) + E_{s_{t+1}}[V^{\pi}(s_{t+1})]$

$$\sum_{t=1}^{n} |S_t, a_t| = r(S_t, a_t) + \sum_{t=1}^{n} |V_{t-1}^{n}|S_{t+1}|$$

 $\approx r(s_t, a_t) + V^{\pi}(s_{t+1})$

Which leads to a simple estimator: $A^{\pi}(s_t, a_t) \approx r(s_t, a_t) + V^{\pi}(s_{t+1}) - V^{\pi}(s_t)$

Thus, estimating is key to computing efficiently

Policy Evaluation via Value Function Fitting

To assess how good a policy is, I evaluate its value function:

$$V_{\frac{\pi}{2}}(\underline{s}_t) = \sum_{\substack{t'=t}}^{T} \mathsf{E}[\underline{r}(\underline{s}_{t'}, a_{t'})|\underline{s}_t]$$

Practically, I approximate this using Monte Carlo returns:

$$V_{(s_t)}^{\pi} \approx \sum_{t=1}^{t} r(s_{t'}, a_{t'})$$

or over N sampled rollouts:

$$V^{\pi}(\underline{s}^{t}) \approx \frac{1}{N} \sum_{i=1}^{N} \frac{1}{t'=t} r(\underline{s}_{i,\underline{t}'}, a_{i,\underline{t}'})$$

Even when using single-sample returns, training a neural network to fit these estimates allows for generalization across similar states, providing a low-variance, learnable baseline to support actor-critic training

1) Monte Carlo Evaluation with Function Approximation

To approximate the value function, I treat it as a supervised learning problem where the target label is the cumulative return from a state. Specifically, the training set is $\{(si,t, yi,t)\}$

Batch Actor-Critic

The actor-critic framework reduces variance in policy gra- dient methods by incorporating a value function (critic). In the batch version (Alg.??), I sample state-action pairs, fit the value function, compute advantages, and update the policy via:

$$\nabla_{\vartheta}(\vartheta) \approx \stackrel{\sim}{\longrightarrow} \nabla_{\vartheta} \log \pi_{\vartheta}(\underline{g}_i|\underline{s}_i) \hat{A}^{\pi}(\underline{s}_i, \underline{g}_i)$$

where $A^{\pi}(s, a) = r(s, a) + V_{\phi}(s') - V_{\phi}(s)$.

Discounting for Infinite Horizons

To prevent unbounded returns in infinite horizon problems, I introduce a discount factor $\gamma \in [0, 1]$. The bootstrapped target while bootstrapped actor-critic estimates reduce variance at the cost of bias. A compromise is to truncate the trajectory and compute n-step returns: becomes yi,t= r(si,t, ai,t) + $\gamma V^{\uparrow} \pi$ (si,t+1). Discounting also affects the policy gradient, where I prefer the form:

$$\nabla J(\vartheta) \approx \frac{1}{N} \int_{N} \nabla \log \pi (a | s) \\ \int_{i=1}^{i=1} \sum_{r=1}^{r=1} I \\ \sum_{i=1}^{r=1} \sum_{r=1}^{r=1} I$$

as it yields lower variance in practice. I incorporate this into the batch actor-critic in Alg.??.GAE ttn=1n t t

Online Actor-Critic

In the online variant (Alg. ??), updates are made per interaction step without storing full trajectories. Using two networks for policy and value function (or a shared network), I update the critic with target $r + \gamma V^{\diamond}(s')$, ompute the advantage, and update the actor with the resulting gradient. For efficiency, parallel simulations are often employed (Fig. ??), though asynchronous setups may result in slightly outdated policies during data collection.

Critics as State-Dependent Baselines

In Monte Carlo policy gradients, the gradient estimate uses the return minus a baseline to reduce variance,

maintaining unbiasedness

$$\nabla \mathscr{Q}_{I}(\vartheta) \simeq \frac{1}{N} \sum_{r=1}^{N} \sum_{t=1}^{T} \nabla_{\vartheta} \log \pi_{\vartheta}(g_{it}|_{Sit}) \qquad (4)$$

$$\sum_{t=1}^{T} \sum_{t=1}^{I} \sum_{t=1}^{N} \sum_{t=1}^{T} \sum_{t=1}^{I} \sum$$

In contrast, actor-critic methods utilize a learned critic to estimate the advantage function, reducing variance at the cost of introducing bias:

$$\nabla \vartheta J(\vartheta) \simeq \frac{1}{N} \sum_{i=1}^{N} \nabla \vartheta \log \pi \vartheta (a_{i,i}|s_{i,i}) \qquad (6)$$
$$r(s_{i,t}, a_{i,t}) + \gamma \hat{V}_{+}^{\pi} (s_{i,t+1}) - \hat{V}_{+}^{\pi} (s_{i,t}) \qquad (7)$$

To better trade off bias and variance, a common approach is to use the critic as a state-dependent baseline, yielding an unbiased yet lower-variance gradient estimator:

$$\bigvee \mathcal{AUU} \cong \frac{1}{N} \sum_{i=1}^{N} \sum_{r=1}^{N} \bigvee_{\vartheta} \log \pi_{\vartheta}(a_{i,\varepsilon}|s_{i,\varepsilon})$$
(8)
$$\sum_{i=1}^{i=1} r=1 \qquad !$$

$$\sum_{r_{i}(s_{i,\varepsilon}, a_{i,\varepsilon})} - \hat{V}_{\phi}^{\pi}(s_{i,\varepsilon}) \qquad (9)$$

Eligibility Traces and n-Step Returns

Monte Carlo and actor-critic methods reflect a bias-variance tradeoff: Monte Carlo returns are unbiased but high-variance, while bootstrapped actor-critic estimates reduce variance at the cost of bias. A compromise is to truncate the trajectory and compute *n*-step returns:

$$\widehat{A_{\pi}}(\underline{s_{t}}, a_{t}) = \underbrace{\sum_{t'=t}^{t'-t} \chi^{t'-t}}_{t' \in t} (\underline{s_{t'}}, a_{t'}) - \widehat{V}_{\phi}^{\pi} (\underline{s_{t}}) + \gamma^{n'} \widehat{V}_{\phi}^{\pi} (\underline{s_{tto}}).$$

Choosing n > 1 often yields more stable learning. General Advantage Estimation (GAE) further refines this by averaging over multiple n-step returns with exponentially decaying weights:

where λ adjusts the tradeoff between bias and variance

6. VALUE FUNCTION METHODS

An Implicit Policy

To bypass policy gradients, I can derive actions directly from the advantage function by selecting the action that maximizes it:

$$\arg_{a_f} \max A^{\pi}(s_t, a_t).$$

This yields an implicit, deterministic policy:

(1, if $a_t = \arg \max_{a_t} A^{\pi}(s_t, a_t)$ $\pi'(a_t | s_t) =$ otherwise

Though the original policy π is not explicitly used, the resulting policy is guaranteed to be at least as good under accurate advantage estimates.

Policy Iteration

1) High-Level Idea

Policy iteration repeatedly improves a policy by evaluating and maximizing the advantage function. Since advantage be expressed in terms of value functions:

 $A^{\pi}(s, a) = r(s, a) + \gamma E [V^{\pi}(s')] - V^{\pi}(s),$

Approximating V $\pi(s)$ suffices to guide the policy update. This method omits gradient steps and instead focuses on alternating

2) **Dynamic Programming**

Assuming discrete state-action spaces and known transition probabilities p(s'|s, a), I can update the value function using the standard dynamic programming (DP) approach:

 $V^{\pi}(s) \leftarrow \mathbb{E}_{g^{\sim}\pi(a|s)} r(s, a) + \gamma \mathbb{E}_{s^{\prime} \sim p(s^{\prime}|s,a)} [V^{\pi}(s^{\prime})]$

If the policy π is deterministic, the expectation over actions disappears, yielding a simplified update:

 $V^{\overline{\alpha}}(s) \leftarrow r(s, \pi(s)) + \gamma E_{s' \sim \rho(s'|s, \pi(s))} [V^{\overline{\alpha}}(s')]$

Alternatively, I can update values directly without main-taining an explicit policy. Since arg maxa $A\pi(s, a) = arg$ maxa $Q\pi(s, a)$ due to $A\pi(s, a) = Q\pi(s, a) - V \pi(s)$, I use this equivalence in value iteration:

Algorithm 3 Policy Iteration (DP)

- while not converged do
 Evaluate V π(s)
- 3. Improve policy $\pi \leftarrow \pi'$

Algorithm 4 Value Iteration (DP)

- 1. while not converged do
- 2. $Q(s, a) \leftarrow r(s, a) + \gamma E[V(s')]$
- 3. $V(s) \leftarrow \max Q(s, a)$

Algorithm 5 Fitted Value Iteration

- 1. while not converged do
- 2. yi \leftarrow maxai (r(si, ai) + $\gamma E[V\phi(s)])$
- 3. Update ϕ to minimize $\Sigma (V\phi(si) yi)2$

Algorithm 6 Fitted Q-Iteration

- 1. while not converged do
- 2. for K times do
- 3. $yi \leftarrow ri + \gamma \max' Q\phi(s', a')$
- 4. Update ϕ to minimize $\Sigma (Q\phi(si, ai) yi)2$

Algorithm 7 Online Q-Iteration

- 1. while interacting with environment do
- 2. Take action ai and observe (si, ai, ri, s')
- 3. $yi \leftarrow ri + \gamma \max' Q\phi(s', a')$
- 4. Update $\phi \leftarrow \phi \alpha \nabla \phi Q \phi(si, ai) (Q \phi(si, ai) yi)$

Fitted Value Iteration

Tabular methods struggle with large state spaces (curse of dimensionality). To address this, I approximate value functions with neural networks. The value function is trained to

$$L(\phi) = \frac{1}{\gamma} V_{\varphi}(s) - \max Q^{\pi}(s, a)^{2}$$

1) Fitted Q-Iteration

To avoid relying on transition dynamics, I learn Q ϕ directly and approximate V $\phi(s) \approx \max Q\phi(s, a)$. This supports off- policy learning and reduces variance, though it lacks convergence guarantees for non-linear approximates. The training minimizes the **Bellman Error**:

 $\epsilon = \frac{1}{2} \mathsf{E}_{(\mathbf{s},\mathbf{g})\sim \mathbf{6}} \underbrace{\mathsf{h}}_{\cdots}(\mathbf{s},a) - \underline{r}(\mathbf{s},a) + \gamma \max_{a} \underbrace{Q_{\mathbf{g}}}_{a} \underbrace{\mathsf{i}_{2}}_{(\mathbf{s}',a')}$

Achieving $\epsilon = 0$ yields the optimal Q-function and policy. However, due to approximation and off-policy sampling, convergence is not guaranteed in general.

2) Online Q-Iteration

Instead of batch learning, I can apply online Q-learning, updating the Q-network with each new sample immediately:

This online, off-policy version enables continual learning with improved sample efficiency but shares the same stability concerns with non-linear function approximates

Value Function Learning Theory

A natural question that arises when exploring value-based methods is whether they converge, and if so, to what. To answer this, I introduce the Bellman backup operator B:

$$BV = \max r_a + \sqrt{T_a}$$

Where ra is the reward vector for action a and Ta is the state transition matrix. The fixed point V * of this operator satisfies:

 $V^* = BV^* = \max r(s, a) + \sqrt{L_1 V^*(s')}$

This fixed point represents the optimal value function. In the tabular case, value iteration converges to V * since B is a contraction under the ∞ norm. In contrast, the non-tabular case introduces challenges. Fit- ted value iteration applies B followed by a projection operator Π onto a function space Ω (e.g., neural networks):

$$IIV = \arg\min_{v \in U} \frac{1}{2} \sum_{(v \cdot (s) - V(s))^2}$$

While both B and Π are contractions in their respective norms, their composition is not, thus convergence is not guaranteed. Fitted Q-iteration uses the same principle: Q $\leftarrow \Pi BQ$, and similarly lacks convergence guarantees due to the compound Operator not being a contraction. Moreover, online Q-iteration updates Q using single samples and a bootstrapped target dependent on the same network, leading to unstable updates.

Replay Buffers and Target Networks

To reduce the correlation in sequential samples, replay buffers B store transitions and allow batch sampling for up However, the bootstrapped target $yi = r + \gamma \max Q\varphi(s, a)$ still depends on the current network parameters φ , leading to biased gradients. To address this, target networks are introduced. A separate parameter set φ' is used to compute targets, typically updated to avoid relying on transition dynamics, I learn $Q\varphi$ directly and approximate $V\varphi(s) \approx \max Q\varphi(s, a)$. This supports off- policy learning and reduces variance, though it lacks convergence guarantees for non-linear approximators. The training minimizes the **Bellman Error**:

$$\phi' \leftarrow \tau \phi' + (1 - \tau) \phi$$

Combining both replay buffers and target networks yields the Deep Q-Network (DQN) algorithm, which significantly improves stability and convergence in practice.

Unifying Q-Learning Variants

The Q-learning family can be generalized by three pro- cesses: data collection, target computation, and parameter updates. In fitted Q-iteration, updates are nested; in online Q-learning, all processes run synchronously; and in DQN, target updates occur at a slower pace than data collection and training. This abstraction helps us better understand and compare these algorithmic variants.

Inaccuracy in Q-Learning

Q-values are not necessarily accurate. The reason lies in the target value. Recall that the target value y is defined as $yj = rj + \gamma \max a' Q\phi'(s', a')$. The max operation in the target is the main problem, because for two random variables X1 and

X2, E $[max(X1, X2)] \ge max$ (E[X1], E[X2]). Therefore, when the next Q-value.

Double Q-Learning

One might notice that $\max a' Q\phi'(s', a') = Q\phi'(s', arg \max a' Q\phi'(s', a'))$. Thus, if I somehow managed to decorrelate the error from the selected action and the error from the Q-function, I could eliminate the erroneous overestimation. To achieve this, I can use two different networks

$$Q_{\phi_A}(s, a) \leftarrow r + \gamma Q_{\phi_B} \quad s', \arg\max_{a} Q_{\phi_A}(s', a')$$

$$Q_{\phi_B}(s, a) \leftarrow r + \gamma Q_{\phi_A} \quad s'_{-} \arg\max_{a} Q_{\phi_A}(s', a')$$

By using the parameters of one network for action selection and the other for value estimation, I decor relate the errors, thereby reducing overestimation bias. In practice, I often use the current and target networks as the two networks. Instead of setting the target as

$$y = r + vQ_{a}(s', arg \max_{a'})$$

I use the current network to select the action and the target network to evaluate its value:

N-Step Return Estimator

In the original definition, the target is $y_i t = r_i t + Q\phi'$ (si,t+1, ai,t+1), which heavily depends on the θ Q-value estimate. When the Q-value estimate is poor, learning stalls. To resolve this, I can use the N -step return trick as in the actor-critic algorithm. The idea is to leverage the bias- variance trade off by limiting the reward accumulation to N steps:

$$y_{it} = \bigvee_{\substack{t'=t}}^{\underline{st}} v^{t'-t} \lim_{t'=t} + \bigvee_{\substack{n \in \mathcal{M} \\ a_{it} \in \mathcal{M}}}^{N} \max_{a_{it} \in \mathcal{M}} Q_{\phi'}(s_{it}, a_{it})$$

However, this introduces an on-policy dependency, as the trajectory of rewards is generated by a specific policy. This limits the ability to fully leverage off-policy data. To mitigate this, I can either:

- Ignore the mismatch (works well in practice),
- Dynamically adapt N to maintain on-policy data,
- Use importance sampling to reweight the returns appro- priately, as described by Munos et al. [2].

Algorithm 8 Deep Deterministic Policy Gradient (DDPG)

- 1. while training do
- 2. Take action ai, observe transition (si, ai, ri, s'), store in replay buffer B
- 3. Sample mini-batch {sj, aj, rj, s' } from B
- 4. Compute target: $yj = rj + \gamma Q\phi'(sj, \mu\theta'(sj))$
- 5. Update critic: $\phi \leftarrow \phi \alpha \nabla \Sigma (Q(s, a) y) 2 \phi \qquad j \qquad j$
- 6. Update actor: $\theta \leftarrow \theta + \beta \Sigma \nabla Q (s, \mu(s))$
- 7. Update target networks

Q-Learning with Continuous Actions

In Q-learning, the implicit policy is defined as:

$$\pi'(a_t|s_t) = \begin{cases} 1, & \text{if } a_t = \arg\max_{a_f} A^{\pi}(s_t, a_t) \\ 0, & \text{otherwise} \end{cases}$$

However, arg max is intractable for continuous action spaces. To resolve this:

Option 1: Sample-based Approximation

Use random sampling over a pre-defined distribution (e.g., uniform) and approximate:

 $\max_{Q(s, a)} \approx \max_{Q(s, a_1), \ldots, Q(s, a_N)}$

This can be improved with techniques like the Cross-Entropy Method (CEM).

Option 2: Structured Q-function

Use a Q-function that is easy to optimize analytically. One such method is the Normalized Advantage Functions (NAF) proposed by Gu et al. [3].

Option 3: Learn an Argmax-er (DDPG)

Train a separate actor network $\mu\theta(s)$ to approximate the action that maximizes the Q-value:

 $\mu_{\theta}(s) \approx \arg \max Q_{\phi}(s, a)$

The optimization becomes:

 $b \leftarrow \arg_{arg} \max Q_{\phi}(s, \mu_{\theta}(s))$

The target becomes: $\mathbf{y}_{i} = \mathbf{r}_{i} \pm \mathbf{y} \mathbf{Q}_{\phi} \left(s_{i'} \mid \mathbf{\mu}_{\theta}(s_{i}) \right) \approx \mathbf{r}_{i} \pm \mathbf{y} \mathbf{Q}_{\phi} \left(s_{i'} \mid \operatorname{arg} \max_{a'} \mathbf{Q}_{\phi} \left(s_{i'} \mid a_{i}' \right) \right)$

4) A Simple ϵ Bound

Assume that $\pi\theta$ is deterministic, i.e., at = $\pi\theta$ (st). From imitation learning, I say $\pi\theta'$ is ϵ -close to $\pi\theta$ if $\pi_{\vartheta}(a_{\vartheta} = \pi_{\vartheta}(s_t) | s_t) \leq \epsilon.$

Then, the new policy's state marginal is given by:

 $p_{\theta}(s_t) = (1 - \epsilon)^t p_{\theta}(s_t) + 1 - (1 - \epsilon)^t p_{\text{mistake}}(s_t).$

I can bound the mismatch in state distributions as:

$$|\underline{p}_{\mathcal{A}}(\underline{s}_{t}) - \underline{p}_{\mathcal{A}}(\underline{s}_{t})| = (1 - (1 - \epsilon)^{t})|\underline{p}_{\text{mistake}}(\underline{s}_{t}) - \underline{p}_{\mathcal{A}}(\underline{s}_{t})|$$

$$\leq 2(1 - (1 - \epsilon)^{t})$$

$$\leq 2\epsilon t$$

This bound is not tight but provides a first-order approximation. Now consider a more general (possibly stochastic) policy $\pi\theta$. I define closeness as:

 $|\pi_{\theta'}(a_t | s_t) - \pi_{\theta}(a_t | s_t)| \le \epsilon,$ ∀st

I use the following lemma: if $|pX(x) - pY(x)| = \epsilon$, then there exists a joint distribution p(x, y) such that $p(x) = \epsilon$. pX (x), p(y) = pY(y), and $p(x = y) = 1 - \epsilon$. Applying this to $\pi\theta$ and $\pi\theta'$, I conclude that the probability they choose different actions is bounded by ϵ . Thus, the same

$$|p_{\alpha}(s_t) - p_{\alpha}(s_t)| \le 2\epsilon t$$

Now, for any function f (st), I can write:

 $\max f(s_i)$

$$\geq \sum_{B_0(s_i)} [f(s_i)] - 2\epsilon t \cdot \max f(s_i).$$

Now apply this to the policy improvement objective:



subject to

$$|\pi_{\vartheta}(a_t | s_t) - \pi_{\vartheta}(a_t | s_t)| \le \epsilon$$
 (12)

A More Convenient Bound - KL Divergence

A tighter and more convenient constraint is provided by KL divergence. Using Pinsker's inequality:

$$|\pi^{\mathfrak{d}}(a \mathfrak{e} | \underline{\mathfrak{s}} \mathfrak{e}) - \pi^{\mathfrak{d}}(a \mathfrak{e} | \mathfrak{s} \mathfrak{e})| \leq \frac{1}{2} D \operatorname{KL}(\pi^{\mathfrak{d}}(\cdot | \mathfrak{s} \mathfrak{e}) || \pi^{\mathfrak{d}}(\cdot | \mathfrak{s} \mathfrak{e}))$$

I define

$$D_{KL}(\underline{p} \| q) = E_{\underline{x} \in \underline{x}} \log \frac{p(x)}{q(x)}$$

Thus, the update becomes:

$$\vartheta' \leftarrow \arg \max_{\sigma} \sum_{\substack{t \in \mathcal{S}_{t} \\ \pi_{\theta}'}} \sum_{\substack{P_{\theta} (s_{t}) \\ T_{\theta}(a_{t} \mid s_{t})}} \sum_{\substack{t \in \mathcal{S}_{t} \\ \pi_{\theta}'}} \sum_{\substack{P_{\theta} (s_{t}) \\ T_{\theta}(a_{t} \mid s_{t})}} \sum_{\substack{t \in \mathcal{S}_{t} \\ T_{\theta}(a_{t} \mid s_{t})}} (13)$$

subject to

υ

ĸ

Can incorporate the constraint using a Lagrangian:

$$L(\vartheta', \lambda) = \sum_{t}^{2} E_{\mathfrak{sline}} E_{\mathfrak{gl}^{\sim}\pi\vartheta} \frac{\pi_{\vartheta'}(a_t \mid \mathfrak{s}_t)}{\pi_{\vartheta}(a_t \mid \mathfrak{s}_t)} \chi^{t} A^{\pi\vartheta}(\mathfrak{s}_t, a_t) -\lambda \left(\underline{D}_{\mathsf{KL}}(\pi_{\vartheta'} \mid \pi_{\vartheta}) - \epsilon \right)$$

I then alternate between:

- Maximizing L with respect to θ' (e.g., via gradient ascent),
- Updating λ via: $\lambda \leftarrow \lambda + \alpha$ (DKL ϵ). This method is known as dual gradient descent.

First-Order Optimization via Taylor Expansion

 $\vartheta' \leftarrow \arg \max_{\vartheta} \nabla_{\vartheta} \overline{A}(\vartheta)^{T}(\vartheta' - \vartheta) \quad \text{s.t.} \quad D_{\mathsf{XL}}(\pi_{\vartheta'} \| \pi_{\vartheta}) \leq \epsilon.$

From policy gradients:

$$\nabla_{\vartheta} \bar{\mathcal{A}}(\vartheta) = \sum_{t} \mathsf{E}_{\mathfrak{s}(t) \mathfrak{s}(t) \mathfrak{s}(t) \mathfrak{s}(t)} \mathcal{V}_{\vartheta}^{t} \nabla_{\vartheta} \log \pi_{\vartheta}(a_{t} \mid \underline{s}_{t}) \underline{\mathcal{A}}^{n\vartheta}(s_{t} \mid a_{t}) \ .$$

Copyright@ REST Publisher

This leads to:

$$\mathcal{J}' \leftarrow \arg \max_{\sigma} \nabla_{\mathcal{J}}(\mathcal{O})^{T} (\mathcal{J} - \mathcal{O}) \quad \text{s.t.} \quad D_{KL}(\pi_{\mathcal{J}} || \pi_{\theta}) \leq \epsilon.$$

Contrast this with the standard gradient ascent step:

 $\vartheta' \leftarrow \arg\max_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} (\vartheta)^{\mathsf{T}} (\vartheta' - \vartheta) \quad \text{s.t.} \quad \left\| \vartheta' - \vartheta \right\|^{\mathsf{L}} \leq \epsilon.$

whose solution is:

$$\mathfrak{g}' \leftarrow \vartheta + \frac{\varepsilon}{\|\nabla \mathfrak{gI}(\vartheta)\|^2} \nabla \mathfrak{gI}(\vartheta).$$

This enforces a spherical constraint in parameter space, rather than in policy space, which is not ideal. I aim reflects the policy mismatch.



FIGURE 8. Dyna: synthetic rollouts from past states

short rollouts provide additional training data. The generalized Dyna algorithm is shown in Algorithm 10. Short synthetic rollouts minimize error accumulation while enhancing sample efficiency by reusing past states for training

Local and Global Models

In the context of LQR, a constrained control optimization problem can be reformulated into an unconstrained one, al- lowing us to minimize cumulative cost terms of the form:

with each xt+1 = f(xt, ut). Solving this via backpropagation requires computing gradients of both the dynamics and cost functions with respect to state and control inputs. For complex systems, exact models are unavailable, so local approximations are used. By linearizing the nonlinear dynamics around a nominal trajectory, I obtain a local model where $f(xt, ut) \approx Atxt + Btut$. These Jacobians, At = df and Bt = df, are estimated using samples from the system.

Local Models

Local models exploit LQR's structure by iteratively updat- ing a linear approximation around the current trajectory. If the true dynamics are stochastic and modeled as p(xt+1|xt, ut) = N (f (xt, ut), Σ), the system can still be approximated linearly within each iteration of iLQR. The resulting feedback control law is given by:

$$u_t = K_t(x_t - \hat{x}_t) + k_t + \hat{u}_{t}$$

where x^t , u^t , Kt, and kt are outputs of the iLQR optimization. To introduce variability and avoid deterministic rollouts, I can inject Gaussian noise into the controller, yielding $p(ut|xt) = N (Kt(xt - x^t) + kt + u^t, \Sigma t)$, where Σt is often set to Q-1. This approach allows for robust exploration during training. short rollouts provide additional training data. The generalized Dyna algorithm is shown in Algorithm 10. Short synthetic rollouts minimize error accumulation while enhancing sample efficiency by reusing past states for training.

LLocal and Global Models

In the context of LQR, a constrained control optimization problem can be reformulated into an unconstrained one, al- lowing us to minimize cumulative cost terms of the form: Dynamics can be modeled through Bayesian linear regres- sion, enabling uncertainty quantification. To maintain policy stability, the updated controller must remain close to previous ones. This is enforced by constraining the divergence be- tween trajectory distributions: $DKL(p(\tau) p(\tau)) \le \epsilon$, ensuring smooth transitions across policy updates.

Guided Policy Search

Guided Policy Search (GPS) bridges local optimal control with global policy learning. The key idea is to use optimized local controllers (e.g., LQR policies) to generate trajectories, which serve as supervised data for training a global policy, typically a neural network. Since a single global policy may not fully capture the behaviors of all local controllers, GPS iteratively refines the local policies while encouraging align- ment with the global policy $\pi\theta$. This is done by modifying the cost to include a regularization term:

 $c_{k+1,i}(x_t, u_t) = c(x_t, u_t) + \lambda_{k+1} \log \pi_{\vartheta}(u_t | x_t),$

which penalizes deviations from the global policy. The overall process is illustrated in Algorithm 11. Over time, both local and global policies are co-optimized to ensure consistency and performance.

Distillation

To efficiently generalize across tasks, reinforcement learning borrows the concept of knowledge distillation from supervised learning [4]. Instead of retaining a collection of specialized models, a single global model is trained to mimic the collective output of an ensemble by learning from their soft probability distributions. The soft targets are produced using:

$$p_i = \frac{\sum_{i \in i} \exp(z_i/T)}{\sum_{i \in i} \exp(z_i/T)},$$

where T is a temperature parameter that smooths the output logits. In RL, this translates to policy distillation, where the global policy πAMN (a|s) is trained to match the behavior of multiple local expert policies πEi (a|s). The loss function is:

$$L = \pi_{\mathcal{E}}(a|s) \log \pi_{AMN}(a|s),$$

which enables the distilled policy to generalize across tasks efficiently. This strategy is foundational for scalable, multi- task RL and has been explored in works like [5], [6]. Similar ideas are also leveraged in Divide-and-Conquer RL, where the role of local LQR controllers is replaced with task-specific neural policies.

Combining Imitation and Reinforcement Learning

Imitation learning offers sample efficiency and stability but is limited by demonstration quality. Reinforcement learning can surpass demonstrations but suffers from exploration chal- lenges. A hybrid approach—pretraining a policy on demon- strations and fine-tuning with reinforcement learning—strikes a balance. However, care must be taken to avoid poor initial trajectories due to distribution shift.

Off-policy Reinforcement Learning

To prevent forgetting demonstrations during training, off- policy RL methods are advantageous, as they allow learning from any data source, including demonstrations reused across iterations. This approach maintains exposure to demonstration data while enabling the learned policy to outperform them, as it is not constrained to imitation. One strategy is to employ off-policy policy gradients, which use importance sampling to adjust for the mismatch between the behavior and target policies:

While policy gradients typically rely on on-policy data, importance sampling enables the use of off-policy data such as demonstrations—by reweighting trajectories. The optimal importance sampling distribution for estimating Ep(x)[f(x)] minimizes variance when $q(x) \propto p(x)|f(x)|$, suggesting that high-reward demonstrations bring us closer to this ideal. To model q(x), I can train a behavioral cloning policy π demo, or if demonstrations come from multiple sources, define a fusion distribution:

$$q(\mathbf{x}) = \frac{1}{M} \sum_{i}^{N} q_i(\mathbf{x})$$

Q-learning with Demonstrations

Since Q-learning is inherently off-policy, it can directly benefit from demonstrations without importance weighting. A practical method is to initialize the replay buffer with demon- stration data and proceed with standard Q-learning updates. This simple augmentation often yields improved performance, especially in sparse reward settings.

Imitation as Auxiliary Loss

Imitation learning maximizes the log-likelihood of expert actions:

 $\sum \log \pi_{\vartheta}(a|s)$

To integrate this with reinforcement learning, a hybrid objective is used:

$$E_{\pi_{\theta}}[r(s, a)] + \lambda \xrightarrow{log \pi_{\theta}(a, s)} \log \pi_{\theta}(a, s)$$

This formulation encourages the agent to imitate expert be- havior while still optimizing for long-term rewards.

Offline reinforcement learning

Unlike traditional RL, which relies on active environment interaction, offline RL aims to learn policies from static datasets—useful in scenarios where data collection is risky or costly, such as autonomous driving. Interestingly, offline RL can sometimes exceed the quality of its dataset by stitching together successful sub-trajectories. For example, Q-learning can converge to an optimal policy even from randomly col- lected data.

REFERENCES

- S. Ross, G. Gordon, and D. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning," in Proceedings of the fourteenth international conference on artificial intelligence and statistics, 2011, pp. 627–635.
- [2]. R. Munos, T. Stepleton, A. Harutyunyan, and M. Bellemare, "Safe and efficient off-policy reinforcement learning," in Advances in Neural Information Processing Systems, 2016, pp. 1054–1062.
- [3]. S. Gu, T. Lillicrap, I. Sutskever, and S. Levine, "Continuous deep q- learning with model-based acceleration," in International Conference on Machine Learning, 2016, pp. 2829–2838.
- [4]. G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," arXiv preprint arXiv:1503.02531, 2015.
- [5]. A. A. Rusu, S. G. Colmenarejo, C. Gulcehre, G. Desjardins, J. Kirk- patrick, R. Pascanu, V. Mnih, K. Kavukcuoglu, and R. Hadsell, "Policy distillation," arXiv preprint arXiv:1511.06295, 2015.
- [6]. E. Parisotto, J. L. Ba, and R. Salakhutdinov, "Actor-mimic: Deep multitask and transfer reinforcement learning," arXiv preprint arXiv:1511.06342, 2015.
- [7]. R. Houthooft, X. Chen, Y. Duan, J. Schulman, F. De Turck, and P. Abbeel, "Vime: Variational information maximizing exploration," in Advances in Neural Information Processing Systems, 2016, pp. 1109–1117.