

Reinforcement Learning: Approximate Dynamic Programming

Decision Making Under Uncertainty, Chapter 10

Christos Dimitrakakis

Chalmers

November 21, 2013

1 Introduction

- Error bounds
- Features

2 Approximate policy iteration

- Estimation building blocks
- The value estimation step
- Policy estimation
- Rollout-based policy iteration methods
- Least Squares Methods

3 Approximate Value Iteration

- Approximate backwards induction
- State aggregation
- Representative states

Definition 1 (u -greedy policy and value function)

$$\pi_u^* \in \arg \max_{\pi} \mathcal{L}_{\pi} u, \quad v_u^* = \mathcal{L} u, \quad (1.1)$$

where $\pi : \mathcal{S} \rightarrow \mathfrak{D}(\mathcal{A})$ maps from states to action distributions.

Parameteric value function estimation

$$\mathcal{V}_{\Theta} = \{v_{\theta} \mid \theta \in \Theta\}, \quad \theta^* \in \arg \min_{\theta \in \Theta} \|v_{\theta} - u\|_{\phi} \quad (1.2)$$

where $\|\cdot\|_{\phi} \triangleq \int_{\mathcal{S}} |\cdot| d\phi$.

Parameteric policy estimation

$$\Pi_{\Theta} = \{\pi_{\theta} \mid \theta \in \Theta\}, \quad \theta^* \in \arg \min_{\theta \in \Theta} \|\pi_{\theta} - \pi_u^*\|_{\phi} \quad (1.3)$$

where $\pi_u^* = \arg \max_{\pi \in \Pi} \mathcal{L}_{\pi} u$

Theorem 2

Consider a finite MDP μ with discount factor $\gamma < 1$ and a vector $u \in \mathcal{V}$ such that $\|u - V_\mu^*\|_\infty = \epsilon$. If π is the u -greedy policy then

$$\|V_\mu^\pi - V_\mu^*\|_\infty \leq \frac{2\gamma\epsilon}{1-\gamma}.$$

In addition, $\exists \epsilon_0 > 0$ s.t. if $\epsilon < \epsilon_0$, then π is optimal.

Feature mapping $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{X}$.

For $\mathcal{X} \subset \mathbb{R}^n$, the feature mapping can be written in vector form:

$$f(s, a) = \begin{bmatrix} f_1(s, a) \\ \dots \\ f_n(s, a) \end{bmatrix} \quad (1.4)$$

Example 3 (Radial Basis Functions)

Let d be a metric on $\mathcal{S} \times \mathcal{A}$ and $\{(s_i, a_i) \mid i = 1, \dots, n\}$. Then we define each element of f as:

$$f_i(s, a) \triangleq \exp \{-d[(s, a), (s_i, a_i)]\}. \quad (1.5)$$

These functions are sometimes called *kernels*.

Example 4 (Tilings)

Let $\mathcal{G} = \{X_1, \dots, X_n\}$ be a **partition** of $\mathcal{S} \times \mathcal{A}$ of size n . Then:

$$f_i(s, a) \triangleq \mathbb{I}\{(s, a) \in X_i\}. \quad (1.6)$$

Approximate policy iteration

Algorithm 1 Generic approximate policy iteration algorithm

input Initial value function v_0 , approximate Bellman operator $\hat{\mathcal{L}}$, approximate value estimator \hat{V} .

for $k = 1, \dots$ **do**

$\pi_k = \arg \min_{\pi \in \hat{\Pi}} \|\hat{\mathcal{L}}_\pi v_{k-1} - \mathcal{L} v_{k-1}\|$ // policy improvement

$v_k = \arg \min_{v \in \hat{V}} \|v - V_\mu^{\pi_k}\|$ // policy evaluation

end for

Theoretical guarantees

Assumption 1

Consider a discounted problem with discount factor γ and iterates v_k, π_k such that:

$$\|v_k - V^{\pi_k}\|_\infty \leq \epsilon, \quad \forall k \quad (2.1)$$

$$\|\mathcal{L}_{\pi_{k+1}} v_k - \mathcal{L} v_k\|_\infty \leq \delta, \quad \forall k \quad (2.2)$$

Theorem 5 ([6], proposition 6.2)

Under Assumption 1

$$\limsup_{k \rightarrow \infty} \|V^{\pi_k} - V^*\|_\infty \leq \frac{\delta + 2\gamma\epsilon}{(1 - \gamma)^2}. \quad (2.3)$$

Lookahead policies

Single-step lookahead

$$\pi_{\mathbf{q}}(a \mid i) > 0 \quad \text{iff } a \in \arg \max_{a' \in \mathcal{A}} q(i, a') \quad (2.4)$$

$$q(i, a) \triangleq r_\mu(i, a) + \gamma \sum_{j \in \mathcal{S}} P_\mu(j \mid i, a) u(j). \quad (2.5)$$

T -step lookahead

$$\pi(i; \mathbf{q}_T) = \arg \max_{a \in \mathcal{A}} q_T(i, a), \quad (2.6)$$

where \mathbf{u}_k is recursively defined as:

$$q_k(i, a) = r_\mu(i, a) + \gamma \sum_{j \in \mathcal{S}} P_\mu(j \mid i, a) \mathbf{u}_{k-1}(j) \quad (2.7)$$

$$\mathbf{u}_k(i) = \max \{q_k(i, a) \mid a \in \mathcal{A}\} \quad (2.8)$$

Rollout policies

Rollout estimate of the q -factor

$$q(i, a) = \frac{1}{K_i} \sum_{k=1}^{K_i} \sum_{t=0}^{T_k-1} r(s_{t,k}, a_{t,k}),$$

where $s_{t,k}, a_{t,k} \sim \mathbb{P}_\mu^\pi(\cdot | s_0 = i, a_0 = a)$, and $T_k \sim \text{Geom}(1 - \gamma)$.

Rollout policy estimation.

Given a set of samples $q(i, a)$ for $i \in \hat{S}$, we estimate

$$\min_{\theta} \|\pi_\theta - \pi_q^*\|_\phi,$$

for some ϕ on \hat{S} .

Generalised linear model using features (or kernel)

Feature mapping $f : \mathcal{S} \rightarrow \mathbb{R}^n$, parameters $\theta \in \mathbb{R}^n$.

$$v_\theta(s) = \sum_{i=1}^n \theta_i f_i(s) \quad (2.9)$$

Fitting a value function.

$$c(\theta) = \sum_{s \in \hat{\mathcal{S}}} c_s(\theta), \quad c_s(\theta) = \phi(s) \|v_\theta(s) - v(s)\|_p^\kappa. \quad (2.10)$$

Example 6

The case $p = 2$, $\kappa = 2$

$$\theta'_j = \theta_j - 2\alpha\phi(s)[v_\theta(s) - v(s)]f_j(s). \quad (2.11)$$

Generalised linear model using features (or kernel).

Feature mapping $f : \mathcal{S} \rightarrow \mathbb{R}^n$, parameters $\theta \in \mathbb{R}^n$.

$$\pi_{\theta}(a | s) = \frac{g(s, a)}{h(s)}, \quad g(s, a) = \sum_{i=1}^n \theta_i f_i(s, a), \quad h(s) = \sum_{b \in \mathcal{A}} g(s, b) \quad (2.12)$$

Fitting a policy through a cost function.

$$c(\theta) = \sum_{s \in \hat{\mathcal{S}}} c_s(\theta), \quad c_s(\theta) = \phi(s) \|\pi_{\theta}(\cdot | s) - \pi(\cdot | s)\|_p^{\kappa}. \quad (2.13)$$

The case $p = 1$, $\kappa = 1$.

$$\theta'_j = \theta_j - \alpha \phi(s) \left(\pi_{\theta}(a | s) \sum_{b \in \mathcal{A}} f_j(s, b) - f_j(s, a) \right). \quad (2.14)$$

Algorithm 2 Rollout Sampling Approximate Policy Iteration.**for** $k = 1, \dots$ **do** Select a set of representative states \hat{S}_k **for** $n = 1, \dots$ **do** Select a state $s_n \in \hat{S}_k$ maximising $U_n(s)$ and perform a rollout. If $\hat{a}^*(s_n)$ is optimal w.p. $1 - \delta$, put s_n in $\hat{S}_k(\delta)$ and remove it from \hat{S}_k . **end for** Calculate $q_k \approx Q^{\pi_k}$ from the rollouts. Train a classifier $\pi_{\theta_{k+1}}$ on the set of states $\hat{S}_k(\delta)$ with actions $\hat{a}^*(s)$.**end for**

Least square value estimation

Projection.

Setting $v = \Phi\theta$ where Φ is a feature matrix and θ is a parameter vector we have

$$\Phi\theta = r + \gamma P_{\mu,\pi}\Phi\theta \quad (2.15)$$

$$\theta = [(I - \gamma P_{\mu,\pi})\Phi]^{-1} r \quad (2.16)$$

Replacing the inverse with the **pseudo-inverse**, with $A = (I - \gamma P_{\mu,\pi})\Phi$

$$\tilde{A}^{-1} \triangleq A^\top (AA^\top)^{-1},$$

Empirical constructions.

Given a set of data points $\{(s_i, a_i, r_i, s'_i) \mid i = 1, \dots, n\}$, which may not be consecutive, we define:

- 1** $r = (r_i)_i$.
- 2** $\Phi_i = f(s_i, a_i)$, $\Phi = (\Phi_i)_i$.
- 3** $P_{\mu,\pi} = P_\mu P_\pi$, $P_{\mu,\pi}(i,j) = \mathbb{I}\{j = i+1\}$

Algorithm 3 LSTDQ - Least Squares Temporal Differences on q -factors

input data $D = \{(s_i, a_i, r_i, s'_i) \mid i = 1, \dots, n\}$, feature mapping f , policy π
 $\theta = (\Phi(\widetilde{I - \gamma P_{\mu, \pi}}))^{-1} r$

Algorithm 4 LSPI - Least Squares Policy Iteration

input data $D = \{(s_i, a_i, r_i, s'_i) \mid i = 1, \dots, n\}$, feature mapping f
Set π_0 arbitrarily.
for $k = 1, \dots$ **do**
 $\theta_k = LSTDQ(D, f, \pi_{k-1})$.
 $\pi_k = \pi_{\Phi \theta_k}^*$.
end for

$$V_t^*(s) = \max_{a \in \mathcal{A}} \{ r(s, a) + \gamma \mathbb{E}_\mu (V_{t+1}^* | s_t = s, a_t = a) \} \quad (3.1)$$

Iterative approximation

$$\hat{V}_t(s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s'} P_\mu(s' | s, a) v_{t+1}(s') \right\} \quad (3.2)$$

$$v_t = \arg \min \left\{ \|v - \hat{V}_t\| \mid v \in \mathcal{V} \right\} \quad (3.3)$$

Online gradient estimation

$$\theta_{t+1} = \theta_t - \alpha_t \nabla_\theta \|v_t - \hat{V}_t\| \quad (3.4)$$

Aggregated estimate.

Let $\mathcal{G} = \{S_0, S_1, \dots, S_n\}$ be a partition of \mathcal{S} , with $S_0 = \emptyset$ and $\theta \in \mathbb{R}^n$ and let $f_k(s_t) = \mathbb{I}\{s_t \in S_k\}$. Then the approximate value function is

$$v(s) = \theta(k), \quad \text{if } s \in S_k, k \neq 0. \quad (3.5)$$

Online gradient estimate.

Consider the case $\|\cdot\| = \|\cdot\|_2^2$. For $s_t \in S_k$:

$$\theta_{t+1}(k) = (1 - \alpha)\theta_t(k) + \alpha \max_{a \in \mathcal{A}} r(s_t, a) + \gamma \sum_j P(j | s_t, a) v_t(s) \quad (3.6)$$

For $s_t \notin S_k$:

$$\theta_{t+1}(k) = \theta(k). \quad (3.7)$$

Representative states approximation.

Let \hat{S} be a set of n representative states and $\theta \in \mathbb{R}^n$ and a feature mapping f :

$$\sum_{i=1}^n f_i(s) = 1, \quad \forall s \in \mathcal{S}.$$

Representative state update.

For $i \in \hat{S}$:

$$\theta_{t+1}(i) = \max_{a \in \mathcal{A}} \left\{ r(i, a) + \gamma \int v_t(s) dP(s | i, a) \right\} \quad (3.8)$$

with

$$v_t(s) = \sum_{i=1}^n f_i(s) \theta_t(i). \quad (3.9)$$

Bellman error methods

$$\min_{\theta} \|v_{\theta} - \mathcal{L}v_{\theta}\| \quad (3.10)$$

Gradient update.

When the norm is

$$\|v_{\theta} - \mathcal{L}v_{\theta}\| = \sum_{s \in \hat{S}} D_{\theta}(s)^2, \quad D_{\theta}(s) = v_{\theta}(s) - \max_{a \in \mathcal{A}} \int_S v_{\theta}(j) dP(j | s, a). \quad (3.11)$$

then the gradient update becomes

$$\theta_{t+1} = \theta_t - \alpha D_{\theta_t}(s_t) \nabla_{\theta} D_{\theta_t}(s_t) \quad (3.12)$$

$$\nabla_{\theta} D_{\theta_t}(s_t) = \nabla_{\theta} v_{\theta_t}(s_t) - \int_S \nabla_{\theta} v_{\theta_t}(j) dP(j | s_t, a_t^*) \quad (3.13)$$

$$a_t^* = \arg \max_{a \in \mathcal{A}} \left\{ r(s_t, a) + \gamma \int_S v_{\theta}(j) dP(j | s_t, a) \right\} \quad (3.14)$$

A litany of approximation algorithms

- Fitted Q-iteration [2].
- Fitted value iteration [16].
- Rollout sampling policy iteration [13]
- State aggregation [19, 4]
- Bellman error minimisation [1, 12, 14]
- Least-squares methods [8, 7, 15].

- [1] A. Antos, C. Szepesvári, and R. Munos. Learning near-optimal policies with bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71(1):89–129, 2008.
- [2] Andrè Antos, Rémi Munos, and Csaba Szepesvari. Fitted q-iteration in continuous action-space mdps. In J.C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20*. MIT Press, Cambridge, MA, 2008.
- [3] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2/3):235–256, 2002.
- [4] A. Bernstein. Adaptive state aggregation for reinforcement learning. Master's thesis, Technion – Israel Institute of Technology, 2007.
- [5] Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control*. Athena Scientific, 2001.
- [6] Dimitri P. Bertsekas and John N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, 1996.
- [7] J.A. Boyan. Technical update: Least-squares temporal difference learning. *Machine Learning*, 49(2):233–246, 2002.
- [8] S.J. Bradtko and A.G. Barto. Linear least-squares algorithms for temporal difference learning. *Machine Learning*, 22(1):33–57, 1996.
- [9] Herman Chernoff. Sequential design of experiments. *Annals of Mathematical Statistics*, 30(3):755–770, 1959.

- [10] Herman Chernoff. Sequential Models for Clinical Trials. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, Vol.4, pages 805–812. Univ. of Calif Press, 1966.
- [11] Morris H. DeGroot. *Optimal Statistical Decisions*. John Wiley & Sons, 1970.
- [12] Christos Dimitrakakis. Monte-carlo utility estimates for bayesian reinforcement learning. In *IEEE 52nd Annual Conference on Decision and Control (CDC 2013)*, 2013. arXiv:1303.2506.
- [13] Christos Dimitrakakis and Michail G. Lagoudakis. Rollout sampling approximate policy iteration. *Machine Learning*, 72(3):157–171, September 2008. Presented at ECML'08.
- [14] Mohammad Ghavamzadeh and Yaakov Engel. Bayesian policy gradient algorithms. In *NIPS 2006*, 2006.
- [15] M.G. Lagoudakis and R. Parr. Least-squares policy iteration. *The Journal of Machine Learning Research*, 4:1107–1149, 2003.
- [16] R. Munos and C. Szepesvári. Finite-time bounds for fitted value iteration. *The Journal of Machine Learning Research*, 9:815–857, 2008.
- [17] Marting L. Puterman. *Markov Decision Processes : Discrete Stochastic Dynamic Programming*. John Wiley & Sons, New Jersey, US, 1994.
- [18] H. Robbins and S. Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, pages 400–407, 1951.
- [19] S. Singh, T. Jaakkola, and M.I. Jordan. Reinforcement learning with soft state aggregation. *Advances in neural information processing systems*, pages 361–368, 1995.

- [20] Niranjan Srinivas, Andreas Krause, Sham Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *ICML 2010*, 2010.
- [21] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.