

2 Stochastic System Models

2.1 Introduction

stochastic system

$$x_{k+1} = f_k(x_k, u_k, w_k), \quad k = 0, 1, 2, \dots$$

discrete state space X

control action space U .

where at time k

$x_k \in X$ is the state,

$u_k \in U$ is the control action.

w_k is the randomness that drives the dynamics, often called plant noise.

x_0 the initial state.

observation of the system is available given by

$$y_k = g_k(x_k, v_k), \quad k = 0, 1, 2, \dots$$

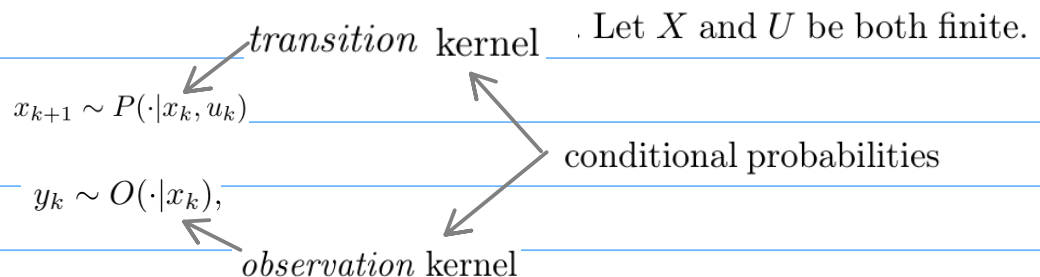
We note that both state and observation equations may change with time, hence the indexing by k . In addition to $\{(f_k, g_k), k = 0, 1, 2, \dots\}$, we must also specify the probability distribution of the random variables $x_0, w_0, w_1, \dots, v_0, v_1, \dots$, and the probability space $(\Omega, \mathcal{F}, \mathbb{P})$ on which they are defined, to complete the specification of a stochastic system. Unless otherwise stated, we will assume these random variables to be *independent*, and for random variables $(w_k, k \geq 0)$ to be *identically distributed* and similarly for random variables $(v_k, k \geq 0)$.

feedback law $\pi = (\pi_0, \pi_1, \dots)$

$$u_k = \pi_k(y_k)$$

Control policies are of two types. When a control action sequence is fixed a priori, say $\bar{u}_0, \bar{u}_1, \dots$, we call it *open-loop control*. The resulting state and observation sequences will be denoted $x_0, \bar{x}_1, \bar{x}_2, \dots$ and $y_0, \bar{y}_1, \bar{y}_2, \dots$ respectively. When the control actions are determined as per a feedback law as in (2.3), we call it *closed-loop control*. A key question is which is better, and whether a particular

2.2 Markov Decision Processes



FACT 2.2 Under any feedback policy π ,

$$P^\pi(X_{k+1} | x^k, u^k) = P^\pi(X_{k+1} | x_k, u_k).$$

FACT 2.3 Under any open-loop policy $\bar{\pi} = (\bar{u}_0, \bar{u}_1, \dots)$,

$$P^\pi(X_{k+m+1}|x^k, \bar{u}^k, \bar{u}_{k+1}^{k+m}) = P^\pi(X_{k+m+1}|x_k, \bar{u}_k, \bar{u}_{k+1}^{k+m}).$$

Markov policy

$$\pi = (\pi_0, \pi_1, \dots)$$

if π_k is a function of x_k alone.

stationary policy

$$\pi_k = \pi$$

i.e., it does not change with time.

finite-horizon Markov Decision Process (MDP).

$$x_{k+1} \sim P(\cdot|x_k, u_k), \quad \text{and} \quad y_k \sim O(\cdot|x_k),$$

with objective function

aka "Reward-to-go"

$$J^\pi = \mathbb{E}\left[\sum_{k=1}^K r(x_k, u_k)\right]$$

infinite-horizon discounted-reward Markov Decision Process (MDP).

$$x_{k+1} \sim P(\cdot|x_k, u_k), \quad \text{and} \quad y_k \sim O(\cdot|x_k),$$

$$J^\pi = \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k r(x_k, u_k)\right],$$

infinite-horizon average-reward Markov Decision Process (MDP).

$$x_{k+1} \sim P(\cdot|x_k, u_k), \quad \text{and} \quad y_k \sim O(\cdot|x_k),$$

$$J^\pi = \liminf_{K \rightarrow \infty} \mathbb{E}\left[\frac{1}{K} \sum_{k=1}^K r(x_k, u_k)\right],$$

Value Function

$$V_0^\pi(i) := r(x_0, u_0) + \mathbb{E}\left[\sum_{k=1}^K r_k(x_k, u_k) | x_0 = i\right]$$

$$\begin{aligned} V_k^\pi(i) &= r_k(i, \pi_k(i)) + \mathbb{E}[V_{k+1}(x_{k+1}) | x_k = i] \\ &= r_k(i, \pi_k(i)) + \sum_j P_{ij}^\pi V_{k+1}(j) \end{aligned}$$

in matrix notation

$$V_k^\pi = r_k^\pi + P_k^\pi V_{k+1}^\pi \quad V_K^\pi = r_K^\pi$$

where $V_k^\pi = (V_k^\pi(1), \dots, V_k^\pi(N))$

FACT 2.6 Under a stationary, Markov policy π ,

$$P^\pi(X_{k+1}, X_{k+2}, \dots | X_k = i) = P^\pi(X_1, X_2, \dots | X_0 = i), \quad \forall i.$$

2.3 Stochastic Linear Systems

stochastic linear system.

$$X = \mathbb{R}^n$$

$$x_{k+1} = Ax_k + Bu_k + w_k,$$

$$U = \mathbb{R}^m$$

$$y_k = Cx_k + v_k,$$

$$Y = \mathbb{R}^p$$

A, B and C are matrices of appropriate dimensions,

at time k , $x_k \in X$ denotes state w_k is plant noise

$u_k \in U$ control action v_k is observation

$y_k \in Y$ observation or sensor noise.

with probability distributions for $x_0, w_0, w_1, \dots, v_0, v_1, \dots$,

ASSUMPTION 2.7 $x_0, w_0, w_1, \dots, v_0, v_1, \dots$ are independent random variables.

ASSUMPTION 2.8 x_0 has Gaussian distribution $\mathcal{N}(\bar{x}_0, \Sigma_0)$. w_k has Gaussian distribution $\mathcal{N}(0, Q)$. v_k has Gaussian distribution $\mathcal{N}(0, R)$.

PROPOSITION 2.9 Suppose A is stable. Then, $\lim_{k \rightarrow \infty} \Sigma_k$ exists, is unique and given by a positive semi-definite matrix, Σ_∞ that is a fixed point of the following equation

$$\Sigma = A\Sigma A^T + Q.$$

3 Finite horizon MDPs

3.1 Introduction

3.2 The Dynamic Programming Algorithm

Algorithm 1 Dynamic Programming

Initial counter: $k = K$.

1. $V_K(x) = r_K(x), \forall x \in X$
2. $V_k(x) = \sup_{u \in U} \{r_k(x, u) + \mathbb{E}[V_{k+1}(f_k(x, u, w_k))]\}, \forall x \in X$
3. $\pi_k(x)$ is maximizer in the above step.
4. while $k > 0, k \leftarrow k - 1$; Goto Step 2.

Output: $V_0(\cdot)$ and $\pi = (\pi_0, \pi_1, \pi_{K-1})$.

LEMMA 3.2 In the backward recursive algorithm (3.4),

$$V_K^\pi(x) = r_K(x), \quad \forall x \in X, \quad (3.4)$$

$$V_k^\pi(x) = r_k(x, u) + \mathbb{E}[V_{k+1}^\pi(f_k(x, u, w_k))], \quad \text{for } k = K - 1, K - 2, \dots, 0.$$

$$V_k^\pi(x_k^\pi) = J_k^\pi(x_k^\pi) \quad a.s. \quad (3.5)$$

THEOREM 3.3 (Optimality of Dynamic Programming) A Markov policy π is optimal if and only if it is a supremizer in Algorithm 1.

Bellman's Principle of Optimality

Optimal Objective Value from k
=
Optimal value (stage reward at k + Optimal Objective Value from k+1)

3.3 The Linear Programming approach to DP

$$\begin{aligned} & \min_{\{v_k\}} \sum_{i=1}^N v_0(i) \\ & \text{s.t. } v_k(i) \geq r_k(i, u) + \sum_{j=1}^N P_{ij}(u) v_{k+1}(j), \quad \forall i, \forall u, \quad k = 0, \dots, K-1, \\ & \quad v_K(i) \geq r_K(i), \quad \forall i. \\ & \quad v_k(i) = \sup_u \{r_k(i, u) + \sum_{j=1}^N P_{ij}(u) v_{k+1}(j)\} \quad V_K(i) = r_K(i) \end{aligned}$$

3.4 Partial Observations and the Belief State

information available at time k .

$$z_k = (y_k, u_{k-1})$$

$$z^k = (z_0, \dots, z_k)$$

belief of the current state

$$p_{k|k}(x_k|z^k) = P(x_k|z^k)$$

and $p_{k+1|k}(x_{k+1}|z^k, u_k) = P(x_{k+1}|z^k, u_k)$

$$\begin{aligned} p_{k+1|k+1}(x_{k+1}|z^{k+1}) &= \frac{P(y_{k+1}|x_{k+1})p_{k+1|k}(x_{k+1}|z^k, u_k)}{\sum_{x_{k+1}} P(y_{k+1}|x_{k+1})p_{k+1|k}(x_{k+1}|z^k, u_k)} \\ &= \frac{P(y_{k+1}|x_{k+1}) \sum_{x_k} P(x_{k+1}|x_k, u_k)p_{k|k}(x_k|z^k)}{\sum_{x_{k+1}} P(y_{k+1}|x_{k+1}) \sum_{x_k} P(x_{k+1}|x_k, u_k)p_{k|k}(x_k|z^k)} \\ &= \Gamma_k[p_{k|k}(\cdot|z^k), u_k, y_{k+1}] \end{aligned}$$

belief (or information) state.

$$\xi_k(z^k) = (p_{k|k}(x_k = i|z^k))_{i=1}^N$$

$$\xi_{k+1}(z^{k+1}) = \Gamma_k[\xi_k(z^k), u_k, y_{k+1}]$$

note

$\xi_k(z^k) \in \Delta(X)$, the space of probability distributions over X

3.5 The DP Algorithm for Partially Observed MDPs

Algorithm 2 Dynamic Programming for POMDPs

Initial counter: $k = K$.

1. $V_K(\xi) = \mathbb{E}[r_K(x)|\xi_K = \xi]$, $\forall \xi \in \Delta(X)$
2. $V_k(\xi) = \sup_{u \in U} \mathbb{E}[r_k(x_k, u) + V_{k+1}(\Gamma_k(\xi_k, y_{k+1}, u))|\xi_k = \xi]$, $\forall \xi \in \Delta(X)$
3. $\pi_k(\xi)$ is a maximizer in the above step.
4. while $k > 0$, $k \leftarrow k - 1$; Goto Step 2.

Output: $V_0(\cdot)$ and $\pi = (\pi_0, \pi_1, \dots, \pi_{K-1})$.

THEOREM 3.4 *Let π^* be a separated policy such that $\pi_k^*(\xi)$ achieves supremum in Step 2 of Algorithm 2. Then, π^* is an optimal policy and $V_k(\xi_k(z^{*,k})) = J_k^*$ a.s.. Conversely, if π^* is a separated optimal policy, then it is a supermizer in Step 2 of Algorithm 2.*

4 Infinite horizon Discounted MDPs

4.1 Introduction

Let us define

$$V^\pi(x_0) = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k r(x_k, u_k) | x_0 \right],$$

then,

$$V^\pi(x_0) = r(x_0, u_0) + \gamma \mathbb{E}[V^\pi(x_1) | x_0],$$

4.2 The Bellman Equation of Optimality

operator $T^\pi : \mathcal{F} \rightarrow \mathcal{F}$

For a stationary Markov policy π

$$[T^\pi v](x) := r(x, \pi(x)) + \gamma \mathbb{E}[v(x_1) | x_0 = x]$$

where \mathcal{F} is the space of all non-negative functions from X to \mathbb{R}

Then,

$$T^\pi V^\pi = V^\pi$$

i.e., V^π is a fixed point of the linear operator T^π .

4.3 An Operator Theoretic View

Bellman operator. $T : \mathcal{F} \rightarrow \mathcal{F}$ such that

$$[Tv](x) := \sup_u \{ r(x, u) + \gamma \mathbb{E}[v(x_1) | x_0 = x] \},$$

\mathcal{F} is the space of all functions $f : X \rightarrow \mathbb{R}^N$

$$TV_\infty = V_\infty \quad \text{i.e., the } V_\infty \text{ is a fixed point}$$

contraction.

$$T : \mathcal{F} \rightarrow \mathcal{F} \quad \text{for any } v_1, v_2 \in \mathcal{F},$$

$$\|Tv_1 - Tv_2\| \leq \gamma \left\| \sum_j (v_1(j) - v_2(j)) P_{ij}(\tilde{u}) \right\| \leq \gamma \left\| \max_j |v_1(j) - v_2(j)| \right\| = \gamma \|v_1 - v_2\|$$

$$\|Tv_1 - Tv_2\| \leq \gamma \|v_1 - v_2\| \quad \sum_j P_{ij}(\tilde{u}) = 1$$

THEOREM 4.1 (Banach Fixed Point Theorem) *Let \mathcal{F} be a complete normed (Banach) space with norm $\|\cdot\|$. Let $T : \mathcal{F} \rightarrow \mathcal{F}$ be a contraction, i.e.,*

$$\|Tv_1 - Tv_2\| \leq \gamma \|v_1 - v_2\|, \quad \forall v_1, v_2 \in \mathcal{F}, \quad 0 < \gamma < 1.$$

Then, (i) there exists a unique fixed point of T , v^ such that $Tv^* = v^*$, and (ii) for any $v_0 \in \mathcal{F}$, the sequence $(v_0, Tv_0, \dots, T^k v_0, \dots)$ has the limit v^* , i.e.,*

$$\lim_{n \rightarrow \infty} \|T^n v_0 - v^*\| = 0.$$

4.4 Dynamic Programming Algorithms

sup-norm,

$$\|v\|_\infty = \max_x |v(x)|.$$

Banach space.

$$(\mathcal{F}, \|\cdot\|_\infty)$$

with $0 \leq r(x, u) \leq \bar{R}.$

V^π bounded by $\bar{V} := \bar{R}/(1 - \gamma)$

$$\mathcal{F} = [0, \bar{V}]^N.$$

$N := \text{cardinality of } X$

The Value Iteration Algorithm

start with any V_0 .

Apply the Bellman operator T to it iteratively, i.e.,

$$V_{k+1}(x) = \sup_u \{r(x, u) + \gamma \sum_{x'} P_{xx'}(u) V_k(x')\}, \quad \forall x.$$

THEOREM 4.2 In the Value Iteration algorithm (4.2), (i)

$$V_n(x) = \sup_{\pi_n} \mathbb{E}^\pi \left[\sum_{k=0}^{n-1} \gamma^k r(x_k, u_k) | x_0 = x \right],$$

and (ii)

$$V^*(x) = \sup_{\pi_n} \mathbb{E}^\pi \left[\sum_{k=0}^{\infty} \gamma^k r(x_k, u_k) | x_0 = x \right].$$

Q-value function,

$$Q^*(x, u) = r(x, u) + \gamma \mathbb{E}[V^*(x') | x, u],$$

V^* is the optimal value function

operator G ,

$$[GQ](x, u) := r(x, u) + \gamma \mathbb{E}[\sup_{u'} Q(x', u') | x, u].$$

note that $V^*(x) = \sup_u Q^*(x, u), \quad \forall x, \quad Q^* = GQ^*,$

The Q-Value Iteration Algorithm

Start with any Q_0 .

Iteratively compute the Q-values,

$$Q_{k+1}(x, u) = r(x, u) + \gamma \sum_{x'} P_{x, x'}(u) \sup_{u'} Q_k(x', u'), \quad \forall x, u.$$

$$\pi^*(x) \in \arg \sup_u Q^*(x, u).$$

The Policy Iteration Algorithm

Start with any (deterministic) policy π_0

policy evaluation

$$\text{solve } T_{\pi_n} V^{\pi_n} = V^{\pi_n}$$

policy improvement

$$\pi_{n+1}(x) \in \arg \sup_u \{ r(x, u) + \gamma \sum_{x'} P_{xx'}(\pi_n(x)) V^{\pi_n}(x') \} \quad \forall x,$$

repeat until $\pi_{n+1} = \pi_n$,

THEOREM 4.3 *The policy iteration algorithm converges to an optimal policy in a finite number of steps. Moreover,*

$$V^{\pi_{n+1}} > V^{\pi_n},$$

at each stage n , with inequality for all states x and strict inequality for at least one x .

Continuous State Spaces *function approximation*

Choose a set of basis functions $\phi_1(x), \dots, \phi_m(x)$

Now approximate V^* as $V^*(x) \approx \sum_{j=1}^m w_j \phi_j(x)$

Pick m points x_1, \dots, x_m

Given V_k , evaluate $V_{k+1}(x_1), \dots, V_{k+1}(x_m)$ as in
the *value iteration* algorithm

Now, find weights $w = (w_1, \dots, w_m)$ that minimize
squared-error, i.e.,

$$\min_w \sum_{i=1}^m (V_{k+1}(x_i) - \sum_{j=1}^m w_j \phi_j(x_i))^2$$

Suppose the solution is \tilde{w} . Then,

$$V_{k+1}(x) \approx \sum_{j=1}^m \tilde{w}_j \phi_j(x)$$

4.6 Empirical Dynamic Programming

Empirical Value Iteration (EVI)

$$\hat{V}_{k+1}(x) = \sup_u \{r(x, u) + \gamma \hat{\mathbb{E}}_n[\hat{V}_k(x')|x, u]\}, \quad \forall x, u.$$

$$\text{with} \quad \hat{\mathbb{E}}_n[V_k(x')|x, u] = \frac{1}{n} \sum_{i=1}^n V_k(x'_i),$$

$x'_i, i = 1, \dots, n$ are samples of the next state from state x with action u .

random operator \hat{T}_n ,

$$\hat{V}_{k+1} = \hat{T}_n \hat{V}_k,$$

$$(\hat{T}_n v)(x) := \sup_u \{r(x, u) + \gamma \hat{\mathbb{E}}_n[v(x')|x, u]\}.$$

By Weak (or Strong) Law of Large numbers, we can expect that $\hat{\mathbb{E}}_n[v(x')|x, u] \rightarrow \mathbb{E}[v(x')|x, u]$ in probability (or almost surely) as $n \rightarrow \infty$. Indeed, we can show that

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(|\hat{T}_n v - T v| > \epsilon \right) = 0, \quad \text{for any } v.$$

4.7 The Linear Programming Approach

linear program.

$$\begin{aligned} & \min \quad \sum_{i=1}^N v(i), \\ & \text{s.t.} \quad v(i) \geq r(i, u) + \gamma \sum_{j=1}^N P_{ij}(u) v(j), \quad \forall i, u. \end{aligned}$$

Note that for each $v(i), i = 1, \dots, N$,

$$v(i) = \sup_u \{r(i, u) + \gamma \sum_{j=1}^N P_{ij}(u) v(j)\},$$

4.8 MDPs with Constraints

constrained MDP problem

$$\begin{aligned} \max_{\pi} \quad & \mathbb{E}[\sum_{k=0}^{\infty} \gamma^k r(x_k, u_k)], \\ \text{s.t.} \quad & \mathbb{E}[\sum_{k=0}^{\infty} \gamma^k c(x_k, u_k)] \leq \bar{C}. \end{aligned}$$

“occupation measure” $\mu(x, u)$

$$\mu(x, u) = \sum_{k=0}^{\infty} \gamma^k P(x_k = x, u_k = u).$$

discounted probability of occupancy in the State-Action space.

CMDP LP

$$\begin{aligned} \max_{\mu} \quad & \sum_{x,u} r(x, u) \mu(x, u), \\ \text{s.t.} \quad & \sum_{x,u} c(x, u) \mu(x, u) \leq \bar{C}, \\ & \sum_u \mu(x, u) = \nu(x) + \gamma \sum_{x'} \sum_{u'} P_{x'|x}(u') \mu(x', u'), \quad \forall x, \\ & \mu(x, u) \geq 0, \quad \forall x, u. \end{aligned}$$

5 Infinite horizon Averaged MDPs

5.1 Introduction

infinite horizon average reward

$$J^\pi := \liminf_{K \rightarrow \infty} \frac{1}{K} \mathbb{E}_\pi \left[\sum_{k=0}^{K-1} r(x_k, u_k) \right]$$

5.2 The Bellman Equation of Optimality

Let us define over a finite horizon K ,

$$V_k^\pi(x) := \mathbb{E}_\pi \left[\sum_{l=0}^k r(x_l, u_l) | x_0 = x \right], \quad k = 0, \dots, K$$

Then, for $u = \pi(x)$,

$$V_K^\pi(x) = \{r(x, u) + \mathbb{E}_\pi[V_{K-1}^\pi(x') | x, u]\}.$$

$$v^\pi(x) = \lim_{K \rightarrow \infty} [V_K^\pi(x) - KJ^\pi], \quad \forall x,$$

$$\lim_{K \rightarrow \infty} \frac{V_K^\pi(x)}{K} = J^\pi,$$

Bellman's equation of optimality for the average reward case:

$$J^* + v^*(x) = \sup_u \{r(x, u) + \mathbb{E}[v^*(x') | x, u]\}, \quad (5.2)$$

v^* is the optimal relative value function

J^* optimal expected average reward

Note that for large horizon K ,

$$V_K^*(x) = KJ^* + v^*(x) + o(1).$$

ASSUMPTION 5.1 The MDP is *unichain*, i.e., for any stationary and Markov policy π , the Markov chain P_π induced by the MDP is irreducible.

Note that an irreducible Markov chain has a single communicating class.

THEOREM 5.2 Under Assumption (5.1), (i) there exists a solution to the Bellman equation (5.2) where J^* is unique and w^* is unique upto an additive constant. (ii) Let π^* denote the policy corresponding to the maximizer in (5.2). Then, π^* is an optimal policy and J^* is the optimal expected average reward. (iii) If π is an optimal policy and $J^* = J^\pi$, then it satisfies (5.2).

LEMMA 5.3 Suppose there exist J^* and v^* that satisfy the average Bellman equation (5.2). Then, J^* is the optimal expected average reward and the supremizing policy π^* is optimal.

LEMMA 5.4 Under Assumption (5.1), there exist solutions (J^*, v^*) to (5.2).

LEMMA 5.5 Suppose π^* is an optimal policy. Then, it satisfies the average Bellman equation (5.2).

Q-relative value function, $q(x, u)$.

$$J^* + q^*(x, u) = r(x, u) + \sum_{x'} P_{xx'}(u) \sup_{u'} q^*(x', u').$$

5.4 DP Algorithms for Average MDPs

Relative Value Iteration Algorithm

Algorithm 3 Relative Value Iteration

Input: $V_0(x) = 0$, $v_0(x) = 0$ and a reference state x_{ref} .

1. $V_{k+1}(x) = \sup_{u \in U} \{r(x, u) + \mathbb{E}[v_k(f_k(x, u, w_k))]\}$, $\forall x \in X$
2. $\pi_{k+1}(x)$ is maximizer in the above step.
3. $v_{k+1}(x) = V_{k+1}(x) - V_{k+1}(x_{ref})$
4. $k \leftarrow k + 1$; Goto Step 2.

Output: $v^*(\cdot)$ and π^* .

relative Q -value iteration (RQVI)

Algorithm 3 Relative Value Iteration

Input: $V_0(x) = 0$, $v_0(x) = 0$ and a reference state x_{ref} .

1. $Q_{k+1}(x, u) = r(x, u) + \mathbb{E}[\sup_{u' \in U} q_k(f_k(x, u, w_k), u')]$, $\forall x \in X$.
2. $\pi_{k+1}(x)$ is maximizer in the above step.
3. $q_{k+1}(x, u) = Q_{k+1}(x, u) - Q_{k+1}(x_{ref}, u_{ref})$, $\forall x, u$.
4. $k \leftarrow k + 1$; Goto Step 2.

Output: $q^*(\cdot)$ and π^* .

The Policy Iteration Algorithm

Start with any (deterministic) policy π_0 .

policy evaluation obtain v^{π_k}

$$J^{\pi_k} + v^{\pi_k}(x) = \{r(x, \pi_k(x)) + \mathbb{E}[v^{\pi_k}(x')|x, \pi_k(x)]\}.$$

setting $v^{\pi_k}(x_{ref}) = 0$ for a reference state x_{ref}

policy improvement

$$\pi_{k+1}(x) \in \arg \sup_u \{r(x, u) + \sum_{x'} P_{xx'}(\pi_k(x)) v^{\pi_k}(x')\} \quad \forall x.$$

repeat until $\pi_{k+1} = \pi_k$,

5.5 Linear Program for Average MDPs

linear program.

$$\begin{aligned} \min_{J^*, w} \quad & J^* \\ \text{s.t.} \quad & J^* + v(i) \geq r(i, u) + \sum_{j=1}^N P_{ij}(u) v(j), \quad \forall i, u. \end{aligned}$$

5.6 Blackwell's Optimality criterion: Average reward as a limit of Discounted reward

THEOREM 5.6 Let V^γ be the optimal value function for the discounted reward MDP with discount factor $\gamma \in (0, 1)$. Then, if

$$|V^\gamma(x') - V^\gamma(x)| < M, \forall x, x',$$

for some $M < \infty$, then there exists a sequence $\{\gamma_l\}$, $0 < \gamma_l < 1$ and $\gamma_l \rightarrow 1$ as $l \rightarrow \infty$ such that limits J^* and v^* exist and satisfy the average reward Bellman equation (5.2).

Bolzano-Weirstrass Theorem,

there exists a sequence $\{\gamma_l\}$,

and the limits

$$J^* =: \lim_{l \rightarrow \infty} (1 - \gamma_l) V^{\gamma_l}(1) \quad \text{and} \\ v^*(x) =: \lim_{l \rightarrow \infty} [V^{\gamma_l}(x) - V^{\gamma_l}(1)] \quad \text{exist.}$$

5.7 MDPs with Constraints

constrained MDP problem of the following kind:

$$\begin{aligned} \max_{\pi} \quad & \lim_{K \rightarrow \infty} \frac{1}{K} \mathbb{E}[\sum_{k=0}^{K-1} r(x_k, u_k)], \\ \text{s.t.} \quad & \lim_{K \rightarrow \infty} \frac{1}{K} \mathbb{E}[\sum_{k=0}^{K-1} c(x_k, u_k)] \leq \bar{C}. \end{aligned}$$

“occupation measure” $\mu(x, u)$,

$$\mu(x, u) := \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=0}^{K-1} P(x_k = x, u_k = u).$$

denoting the initial state distribution by ν ,

CMDP LP average reward

$$\begin{aligned} \max_{\mu} \quad & \sum_{x, u} r(x, u) \mu(x, u), \\ \text{s.t.} \quad & \sum_{x, u} c(x, u) \mu(x, u) \leq \bar{C}, \\ & \sum_u \mu(x, u) = \sum_u \sum_{x', u'} \sum_{u'} P_{x'x}(u') \mu(x', u'), \quad \forall x, \\ & \mu(x, u) \geq 0, \quad \forall x, u. \end{aligned}$$

5.8 An EDP Algorithm

Empirical Relative Value Iteration (ERVI)

$$\hat{V}_{k+1}(x) = \sup_u \{r(x, u) + \hat{\mathbb{E}}_n[\hat{v}_k(x')|x, u]\}, \quad \forall x, u. \quad \hat{v}_{k+1}(x) = \hat{V}_{k+1}(x) - \hat{V}_{k+1}(x_{ref}),$$

x_{ref} is a reference state (chosen arbitrarily)

with

$$\hat{\mathbb{E}}_n[v_k(x')|x, u] = \frac{1}{n} \sum_{i=1}^n v_k(x'_i)$$

$x'_i, i = 1, \dots, n$ are samples of the next state from state x with action u .

random operator \hat{T}_n

$$\hat{v}_{k+1} = \hat{T}_n \hat{v}_k,$$

$$(\hat{T}_n v)(x) := \sup_u \{r(x, u) + \gamma \hat{\mathbb{E}}_n[v(x')|x, u] - \sup_u \{r(x_{ref}, u) + \gamma \hat{\mathbb{E}}_n[v(x')|x_{ref}, u]\}.$$

