# SEMI-MARKOV DECISION PROCESSES

M. Baykal-Gűrsoy


Department of Industrial and Systems Engineering
Rutgers University
Piscataway, New Jersey
Email: gursoy@rci.rutgers.edu

### Abstract

Considered are infinite horizon semi-Markov decision processes (SMDPs) with finite state and action spaces. Total expected discounted reward and long-run average expected reward optimality criteria are reviewed. Solution methodology for each criterion is given, constraints and variance sensitivity are also discussed.

## 1 Introduction

Semi-Markov decision processes (SMDPs) are used in modeling stochastic control problems arising in Markovian dynamic systems where the sojourn time in each state is a general continuous random variable. They are powerful, natural tools for the optimization of queues [20, 44, 41, 18, 42, 43, 21], production scheduling [35, 31, 2, 3], reliability/maintenance [22].

For example, in a machine replacement problem with deteriorating performance over time, a decision maker, after observing the current state of the machine, decides whether to continue its usage, or initiate a maintenance (preventive or corrective) repair, or replace the machine. There are a reward or cost structure associated with the states and decisions, and an information pattern available to the decision maker. This decision depends on a performance measure over the planning horizon which is either finite or infinite, such as total expected discounted or long-run average expected reward/cost with or without external constraints, and variance penalized average reward.

SMDPs are based on semi-Markov processes (SMPs) [9] [Semi-Markov Processes], that include renewal processes [Definition and Examples of Renewal Processes] and continuous-time Markov chains (CTMCs) [Definition and Examples of CTMCs] as special cases. In a semi-Markov process similar to Markov chains (DTMCs) [Definition and Examples of DTMCs], state changes occur according to the Markov property, i.e., states in the future do not depend on the states in the past given the present. However, the sojourn time in a state is a continous random variable with distribution depending on that state and the next state, a Markov chain is a SMP in which the sojourn times are discrete (geometric) random variables independent of the next state; a continuous-time Markov chain is an SMP with exponentially distributed sojourn times; and a renewal process is an SMP with a single state. Semi-Markov decision processes, first introduced by Jewell [23] and De Cani [8], are also called as Markov renewal programs [10, 12, 19, 30, 37].

This article is organized as follows. The next section introduces basic definitions and notations. Various performance criteria are presented in Section 3 and their solution methodologies are described in Sections 4-6.

## 2  Basic Definitions

We consider time-homogeneous, finite state and finite action SMDPs, and give references for the more general cases. Let $\{X_m, m \geq 0\}$ denote the state process, which takes values in a finite state space $\mathcal{S}$. We also use $\{X_m, m \in \mathcal{N}\}$ to denote the state process with $\mathcal{N}$ representing the set of nonnegative integers. At each epoch $m$ the decision maker chooses an action $A_m$ from a finite action space $\mathcal{A}$. The sojourn time between the $(m-1)$-st and the $(m)$-th epochs is a random variable and denoted by $\Upsilon_m$. The underlying sample-space $\Omega = \{\mathcal{S} \times \mathcal{A} \times (0,\infty)\}^\infty$ consists of all possible realizations of states, actions and the transition times. Throughout, the sample space will be equipped with the $\sigma$-algebra generated by the random variables $\{X_m, A_m, \Upsilon_{m+1}; \; m \geq 0\}$. The initial state is assumed to be fixed and given. Note that we will suppress the dependence on the initial state unless given otherwise. Denote $P_{xay}$, $x \in \mathcal{S}$, $a \in \mathcal{A}$, $y \in \mathcal{S}$, for the law of motion of the process, i.e., for all policies $\boldsymbol{u}$ and all epochs $m$

$$P_{\boldsymbol{u}}\{X_{m+1} = y | X_0, A_0, \Upsilon_1, \ldots, X_m = x, A_m = a\} = P_{xay}.$$

Also conditioned on the event that the next state is $y$, $\Upsilon_{m+1}$ has the distribution function $F_{xay}(.)$, i.e.,

$$P_{\boldsymbol{u}}\{\Upsilon_{m+1} \leq t | X_0, A_0, \Upsilon_1, \ldots, X_m = x, A_m = a, X_{m+1} = y\} = F_{xay}(t).$$

Assume that $F_{xay}(0) < 1$.

The process $\{S_t, B_t : t \geq 0\}$ where $S_t$ is the state of the process at time $t$, and $B_t$ is the action taken at time $t$, is referred to as the Semi-Markov Decision Process. Let $T_n = \sum_{m=1}^n \Upsilon_m$, i.e., denote the time of $n$-th transition. For $t \in [T_m, T_{m+1})$, clearly

$$S_t = X_m, \quad B_t = A_m.$$

### 2.1  Policy Types

A *decision rule* $\boldsymbol{u}^m$ at epoch $m$ is a vector consisting of probabilities assigned to each available action. A decision rule may depend on all of the previous states, actions, transition times and the present state. Let $u_a^m$ denote the $a$-th component of $\boldsymbol{u}^m$. Thus, it is the conditional probability of choosing action $a$ at the $m$-th epoch, i.e.,

$$P_{\boldsymbol{u}}\{A_m = a | X_0 = x_0, A_0 = a_0, \Upsilon_1 = \tau_1, \ldots, X_m = x\} = u_a^m(x_0, a_0, \tau_1, \ldots, x).$$

A *policy* is an infinite sequence of decision rules $\boldsymbol{u} = \{\boldsymbol{u}^0, \boldsymbol{u}^1, \boldsymbol{u}^2, \ldots\}$.

Policy $\boldsymbol{u}$ is called Markov policy if $\boldsymbol{u}^m$ at epoch $m$ depends only on the current state not the past history, i.e.,

$$\boldsymbol{u}^m(x) = \boldsymbol{u}_a^m(x_0, a_0, \tau_1, \ldots, x).$$

A policy is called *stationary* if the decision rule at each epoch is the same and it depends only on the present state of the process, $\boldsymbol{u} = \{\boldsymbol{u}, \boldsymbol{u}, \boldsymbol{u}, \ldots\}$; denote $f_{xa}$ for the probability of choosing action $a$ when in state $x$. A stationary policy is said to be *pure* if for each $x \in \mathcal{S}$ there is only one action $a \in \mathcal{A}$ such

that $f_{xa} = 1$. Let $U$, $M$, $F$ and $G$ denote the set of all policies, Markov policies, stationary policies and pure policies, respectively. Clearly, $G \subset F \subset M \subset U$.

Under a stationary policy $\boldsymbol{f}$ the state process $\{S_t : t \geq 0\}$ is a semi-Markov process, while the process $\{X_m : m \in \mathcal{N}\}$ is the embedded Markov chain with transition probabilities

$$P_{xy}(\boldsymbol{f}) = \sum_{a \in \mathcal{A}} P_{xay} f_{xa}.$$

Clearly, the process $\{S_t, B_t : t \geq 0\}$ is also a semi-Markov process under a stationary policy $\boldsymbol{f}$ with the embedded Markov chain $\{X_m, A_m : m \in \mathcal{N}\}$.

## 2.2  Chain Structure

Under a stationary policy $\boldsymbol{f}$, state $x$ is *recurrent* if and only if $x$ is recurrent in the embedded Markov chain; similarly, $x$ is *transient* if and only if $x$ is transient for the embedded Markov chain. A semi-Markov decision process is said to be *unichain(multi-chain)* if the embedded Markov chain for each pure policy is unichain (multi-chain), i.e., if the transition matrix $P(\boldsymbol{g})$ has at most one (more than one) recurrent class plus (a perhaps empty) set of transient states for all pure policies $\boldsymbol{g}$. It is called *irreducible* if $P(\boldsymbol{g})$ is irreducible under all pure policies $\boldsymbol{g}$. Similarly, an SMDP is said to be *communicating* if $P(\boldsymbol{f})$ is irreducible for all stationary policies that satisfy $f_{xa} > 0$, for all $x \in \mathcal{S}$, $a \in \mathcal{A}$.

Let $\tau(x, a)$ define the expected sojourn time given that the state is $x$ and the action $a$ is chosen just before a transition, i.e.,

$$
\begin{aligned}
\tau(x, a) \quad &\triangleq \quad E_{\boldsymbol{u}}[\Upsilon_m | X_{m-1} = x, A_{m-1} = a] \\
&= \int_0^\infty \sum_{y \in \mathcal{S}} P_{\boldsymbol{u}} \{X_m = y, \Upsilon_m > t | X_{m-1} = x, A_{m-1} = a\} dt \\
&= \int_0^\infty [1 - \sum_{y \in \mathcal{S}} P_{xay} F_{xay}(t)] dt.
\end{aligned}
$$

Let $W_t(x, a)$ denote the random variables representing the state-action intensities,

$$W_t(x, a) \triangleq \frac{1}{t} \int_0^t \mathbf{1}\{(S_s, B_s) = (x, a)\} \, ds,$$

where $\mathbf{1}\{.\}$ denotes the indicator function. Let $U_0$ denote the class of all policies $\boldsymbol{u}$ such that $\{W_t(x, a); t \geq 0\}$ converges. Thus, for $\boldsymbol{u} \in U_0$, there exist random variables $\{W(x, a)\}$ such that

$$\lim_{t \to \infty} W_t(x, a) = W(x, a).$$

Let $U_1$ be the class of all policies $\boldsymbol{u}$ such that the expected state-action intensities $\{E_{\boldsymbol{u}}[W_t(x, a)]; t \geq 0\}$ converge for all $x$ and $a$. For $\boldsymbol{u} \in U_1$ denote

$$w_{\boldsymbol{u}}(x, a) = \lim_{t \to \infty} E_{\boldsymbol{u}}[W_t(x, a)].$$

From Lebesgue's Dominated Convergence Theorem $U_0 \subset U_1$.

3

A well-known result from renewal theory (see Çınlar [9]) is that if $\{Y_t = (S_t, B_t) : t \geq 0\}$ is a homogeneous semi-Markov process, and if the embedded Markov chain $\{X_m, m \in \mathcal{N}\}$ is unichain then, the proportion of time spent in state $y$, i.e.,

$$\lim_{t \to \infty} \frac{1}{t} \int_0^t \mathbf{1}\{Y_s = y\} \, ds,$$

exists. Since under a stationary policy $\boldsymbol{f}$ the process $\{Y_t = (S_t, B_t) : t \geq 0\}$ is a homogeneous semi-Markov process, if the embedded Markov decision process is unichain then the limit of $W_t(x, a)$ as $t$ goes to infinity exists and the proportion of time spent in state $x$ when action $a$ is applied is given as

$$W(x, a) = \lim_{t \to \infty} W_t(x, a) = \frac{\tau(x, a) Z(x, a)}{\sum_{x,a} \tau(x, a) Z(x, a)},$$

where $Z(x, a)$ denotes the associated state action frequencies. Let $\{z_{\boldsymbol{f}}(x, a); x \in \mathcal{S}, a \in \mathcal{A}\}$ denote the expected state-action frequencies, i.e.,

$$z_{\boldsymbol{f}}(x, a) = \lim_{n \to \infty} E_{\boldsymbol{f}} \frac{1}{n} \sum_{m=1}^n \mathbf{1}\{X_{m-1} = x, A_{m-1} = a\} = \pi_x(\boldsymbol{f}) f_{xa}$$

where $\pi_x(\boldsymbol{f})$ is the steady-state distribution of the embedded Markov chain $P(\boldsymbol{f})$.

The long-run average number of transitions into state $x$ when action $a$ is applied per unit time is,

$$v_{\boldsymbol{f}}(x, a) = \frac{\pi_x(\boldsymbol{f}) f_{xa}}{\sum_{x,a} \tau(x, a) \pi_x(\boldsymbol{f}) f_{xa}} = \frac{z_{\boldsymbol{f}}(x, a)}{\sum_{x,a} \tau(x, a) z_{\boldsymbol{f}}(x, a)}. \tag{1}$$

This gives $w_{\boldsymbol{f}}(x, a) = \tau(x, a) v_{\boldsymbol{f}}(x, a)$.

## 2.3   Reward Structure

Let $R_t$ be the reward function at time $t$. $R_t$ can be an impulse function corresponding to the reward earned immediately at a transition epoch and/or it can be a step function between transition epochs corresponding to the rate of reward as described below. The decision maker earns an immediate reward $R(X_m, A_m)$ and a reward with rate $r(X_m, A_m)$ until the $(m+1)$-th epoch, i.e.,

$$R_t = \begin{cases} R(X_m, A_m), & \text{if } t = T_m, \\[2mm] r(X_m, A_m), & \text{if } t \in [T_m, T_{m+1}). \end{cases} \tag{2}$$

Thus,

$$R_{m+1} = R(X_m, A_m) + r(X_m, A_m) \Upsilon_{m+1},$$

is the reward earned during the $(m+1)$-th transition [6, 40, 45].

Similarly, there is an immediate cost $C(X_m, A_m)$ and a cost with rate $c(X_m, A_m)$ with

$$C_{m+1} = C(X_m, A_m) + c(X_m, A_m) \Upsilon_{m+1}.$$

Hence, at any epoch if the process is in state $x \in \mathcal{S}$ and action $a \in \mathcal{A}$ is chosen, then the reward earned during this epoch is represented by $\bar{r}(x, a) \stackrel{\triangle}{=} R(x, a) + r(x, a)\tau(x, a)$. Similarly, the cost during this epoch is represented by $\bar{c}(x, a) \stackrel{\triangle}{=} C(x, a) + c(x, a)\tau(x, a)$.

4

**Example:** Consider the *machine replacement problem* mentioned in the Introduction with states (1) machine is in good condition, (2) machine has some minor problems, (3) machine is down and needs to be replaced. Time to failure of the machine follows a Weibull distribution with scale parameter equal to 8,000 hours and shape parameter equal to 4. The failure is minor with probability .95 and major requiring replacement of the machine with probability 0.05. Life time of a machine with minor problems follows a Weibull distribution with scale parameter equal to 20,000 hours and shape parameter equal to 4. However, a machine with minor problems could be maintenance repaired to it as good as new. Maintenance repair takes Weibull distributed amount of time with scale parameter equal to 8 hours and shape parameter equal to 0.5. On the other hand, the machine replacement time is normally distributed with mean 72 hours and variance of 8 hours. Running a fully working machine earns $100/hr$, and a machine with minor problem earns $75/hr$ profit. It costs $40/hr$ to repair and $10,000$ to replace a machine. Note that there is no control action available in states 1 and 3. In state 1 the decision maker needs to "wait" and in state 3 s/he needs to order a new machine. Let us denote these action as action 1. In state 2, there are two possible actions to choose: "wait" action denoted as action 1, and "initiate repair" denoted as action 2. Parameters of this model are:

$$P_{113} = 0.95, \quad P_{113} = 0.05, \quad P_{213} = 1, \quad P_{221} = 1, \quad P_{311} = 1,$$
$$\tau(1,1) = 7,252, \quad \tau(2,1) = 1,813, \quad \tau(2,2) = 48, \quad \tau(3,1) = 72,$$
$$\bar{r}(1,1) = 725,200, \quad \bar{r}(2,1) = 135,975, \quad \bar{r}(2,2) = -1,920, \quad \bar{r}(3,1) = -10,000.$$

The last two reward values correspond to the incurred costs under repair and replacement, respectively.

## 3    Performance Measures

We will focus on the optimality criteria over the infinite horizon, since some general results could be obtained for these models. We will first consider finding a policy $\boldsymbol{u}$ that will maximize the *total discounted reward* defined as

$$\phi_\alpha(\boldsymbol{u}) \triangleq E_{\boldsymbol{u}}[\int_0^\infty e^{-\alpha s} R_s \, ds]. \tag{3}$$

where $\alpha$ represents the discount factor[23, 39, 17, 26, 28]. Discounted reward optimality criterion is easier to analyze and understand than the average reward criterion, since the results for these models hold regardless of the chain structure of the embedded Markov chain. In fact, the existence of this integral is immediate under finite rewards. In addition, discounting lands itself naturally to economic problems in which the present value of future earnings is discounted as a function of the interest rate. Another interpretation of these models implies the importance of the initial decisions.

The great majority of the literature, on the other hand, is concerned with the *long-run average expected reward* criterion with

$$\phi_1(\boldsymbol{u}) \triangleq \liminf_{t \to \infty} \frac{1}{t} E_{\boldsymbol{u}}[\int_0^t R_s \, ds], \tag{4}$$

$\phi_1$ denoting the *long-run average expected reward* [23, 19, 11, 12, 27, 38, 45]. The following alternative to $\phi_1$ is given by Jewell [24], Ross [34, 33], and Mine and Osaki [30] as

$$\phi_2(\boldsymbol{u}) \triangleq \liminf_{n \to \infty} \frac{E_{\boldsymbol{u}}[\sum_{m=1}^n R_m]}{E_{\boldsymbol{u}}[T_n]}, \tag{5}$$

referred to as the *ratio-average reward* [16]. The performance measure $\phi_2$ is also used by other researchers (see e.g., [7, 15, 13, 14, 23, 21, 32]).

Let

$$\phi_\alpha^* = \sup_{\boldsymbol{u} \in U} \phi_\alpha(\boldsymbol{u}), \quad \phi_1^* = \sup_{\boldsymbol{u} \in U} \phi_1(\boldsymbol{u}), \quad \phi_2^* = \sup_{\boldsymbol{u} \in U} \phi_2(\boldsymbol{u}).$$

A policy $\boldsymbol{u}$ is optimal for $\phi_\alpha(\cdot)$ if $\phi_\alpha(\boldsymbol{u}) = \phi_\alpha^*$. For a fixed $\epsilon > 0$, a policy $\boldsymbol{u}$ is $\epsilon$-optimal for $\phi_\alpha(\cdot)$ if $\phi_\alpha(\boldsymbol{u}) > \phi_\alpha^* - \epsilon$. Optimality and $\epsilon$-optimal for $\phi_1(\cdot)$, $\phi_2(\cdot)$ and the other performance measures we will consider in this article are defined analogously.

The following *expected time-average reward* criterion has been considered recently by Baykal-Gürsoy and Gürsoy [1, 4]( also see [5])

$$\psi(\boldsymbol{u}) \triangleq E_{\boldsymbol{u}}[\liminf_{t \to \infty} \frac{1}{t} \int_0^t R_s \, ds] \tag{6}$$

subject to the sample path constraint,

$$P_{\boldsymbol{u}}\{\limsup_{t \to \infty} \frac{1}{t} \int_0^t C_s \, ds \leq \gamma\} = 1. \tag{7}$$

This constraint requires that the long-run average costs on almost all sample paths should be bounded by $\gamma$.

More generally, they investigate the following *expected time-average variability*

$$\nu(\boldsymbol{u}) \triangleq E_{\boldsymbol{u}}[\liminf_{t \to \infty} \frac{1}{t} \int_0^t h(R_s, \frac{1}{t} \int_0^t R_q \, dq) \, ds], \tag{8}$$

where $h(.,.)$ is a continuous function of the current reward at time $s$ and the average reward over an interval that includes time $s$. By letting $\nu^* = \sup_{\boldsymbol{u} \in U} \nu(\boldsymbol{u})$, the optimality and $\epsilon$-optimality for $\nu(.)$ are analogously defined.

# 4    Discounted Reward Criterion

Discounted reward can be rewritten as:

$$
\begin{aligned}
\phi_\alpha(\boldsymbol{u}) \quad &= E_{\boldsymbol{u}}[\sum_{m=0}^\infty e^{-\alpha T_m}(R(X_m, A_m) + \frac{r(X_m, A_m)}{\alpha}(1 - e^{-\alpha \Upsilon_m}))] \\
&= \sum_{m=0}^\infty \sum_{x,a} \int_0^\infty e^{-\alpha t}[R(x,a) + \frac{r(x,a)}{\alpha}(1 - \sum_y P_{xay} \int_0^\infty e^{-\alpha \tau} dF_{xay}(\tau))]P_{\boldsymbol{u}}\{X_m = x, A_m = a, T_m \leq t\}.
\end{aligned}
$$

The terms inside the second integral could be recognized as the Laplace transform of the density function $f_{xay}(\cdot)$ and will be denoted as $\tilde{f}_{xay}(\alpha)$.

The optimal discounted reward vector is represented by $\phi_\alpha^{*x}$ for each initial state $x$, and it can be shown that it satisfies the optimality equation for all $x \in \mathcal{S}$ :

$$
\begin{aligned}
\phi_\alpha^x \quad &= \max_a\{[R(x,a) + \frac{r(x,a)}{\alpha}(1 - \sum_j P_{xaj}\tilde{f}_{xaj})] + \sum_y P_{xay}\tilde{f}_{xay}(\alpha)\phi_\alpha^y\} \\
&= \max_a\{r^\alpha(x,a) + \sum_y P_{xay}^\alpha \phi_\alpha^y\}. \tag{9}
\end{aligned}
$$

Second equality is obtained from the first by denoting the terms inside the square bracket as $r^\alpha(x,a)$ and writing $P_{xay}\tilde{f}_{xay}(\alpha)$ as $P^\alpha_{xay}$. Note that the second equality is similar to the one obtained for the Markov Decision Processes(MDPs) [The Total Expected Discounted Reward MDPs: Existence of Optimal Policies]. Thus, discounted SMDPs can be reduced to discounted MDPs by using these transformations. Since $P^\alpha_{xay} < 1$, the right hand side of the optimality Equation (9) is a contraction mapping and the next theorem is immediate.

**Theorem 1** *For SMDPs under the discounted reward criterion:*

(i) *There exists a unique solution to the optimality equation 9 and it is equal to $\phi^*_\alpha$.*

(ii) *There exists an optimal pure policy $\boldsymbol{g}^*$ given by, $\phi_\alpha(\boldsymbol{g}^*) = \phi^*_\alpha = (I - P^\alpha_{xay})^{-1} r^\alpha(\boldsymbol{g}^*)$ where $r^\alpha(\boldsymbol{g}^*)$ denotes the single-period discounted reward earned under policy $\boldsymbol{g}^*$.*

This optimal pure policy could be obtained using the policy iteration [The Total Expected Discounted Reward MDPs: Policy Iteration], value iteration [The Total Expected Discounted Reward MDPs: Value Iteration] or linear programming algorithms [21, 42, 43, 32] [Linear Programming Formulations of MDPs]. The linear programming(LP) algorithm is discussed next. Consider the following LP with given numbers $\beta_x > 0$ for $x \in \mathcal{S}$ .

$$max \quad \sum_{x \in \mathcal{S},a \in \mathcal{A}} r^\alpha(x,a)z(x,a)$$

$$s.t. \quad \sum_{x \in \mathcal{S},a \in \mathcal{A}} (\delta_{xy} - P^\alpha_{xay})z(x,a) = \beta_y, \quad y \in \mathcal{S}$$

$$z(x,a) \geq 0, \quad x \in \mathcal{S}, \ a \in \mathcal{A} .$$

Let $\boldsymbol{z}^*$ be an optimum solution of the above LP. Clearly, any extreme point of this LP has $|\mathcal{S}|$ number of basic variables where $|\cdot|$ denotes the number of elements in a given set. Thus, $z^*(x,a)$ is positive only for one action $a$. The optimum pure policy $\boldsymbol{g}^*$ is then obtained by assigning $\boldsymbol{g}^*_x$ in such a way that $z^*(x, \boldsymbol{g}^*_x) > 0$. Constraints defined in a similar fashion,

$$E\boldsymbol{u}\Big[\int_0^\infty e^{-\alpha s}C_s\, ds\Big] < \gamma,$$

could be included into the LP as

$$\sum_{x \in \mathcal{S},a \in \mathcal{A}} c^\alpha(x,a)z(x,a) < \gamma,$$

with $c^\alpha(x,a) = C(x,a) + \frac{c(x,a)}{\alpha}(1 - \sum_y P_{xay}\tilde{f}_{xay})$. Since every new constraint will increase the number of basic variables, the optimum policy will no longer be pure but randomized stationary [17] [Constrained MDPs].

For the countable state case we need the assumption,

**Assumption 1** *There exists $\delta > 0$ and $\varepsilon > 0$, such that*

$$F_{xay}(\delta) \leq 1 - \varepsilon \quad for \ all \ x \ and \ y \in \mathcal{S} \ \ and \ a \in \mathcal{A} ,$$

together with $|r^\alpha(x,a)| \leq M < \infty$ to ensure the existence of an optimal pure policy. Additional conditions are required for SMDPs with Borel state and action spaces, and unbounded rewards [32, 17, 40].

# 5 Average Reward Criterion

*Average* or *ratio-average* expected reward criterion is applied to systems in which the system dynamics is not slow enough to warrant discounting. This criterion is more difficult to analyze since the existence of the optimal stationary policy depends on the chain structure. Under the condition that the SMDP is irreducible, $\phi_1(\boldsymbol{f}) = \phi_2(\boldsymbol{f})$ for every stationary policy $\boldsymbol{f}$ [34, 30]. However, this may not hold even for unichain SMDPs [16]. While $\phi_1$ is clearly the more appealing criterion, it is easier to write the optimality equations when establishing the existence of an optimal pure policy under criterion $\phi_2$ [38, 45, 36]. On the other hand, for finite state and finite action SMDPs there exists an optimal pure policy under $\phi_1$ [12, 38, 36], while such an optimal policy may not exists under $\phi_2$ in a general multichain SMDP [25]. Jianyong and Xiaobo [25] investigate average reward SMDPs focusing on $\phi_2$ and using a data-transformation method [37]. They show that the optimal pure policy exists in some special cases such as the unichain case and the weakly communicating case.

The optimal pure policy for the average expected reward criterion in multi-chain SMDPs is obtained from the optimal solution of the following LP [26] under the assumption on the sojourn times.

$$
\begin{aligned}
max \quad & \sum_{x \in \mathcal{S}, a \in \mathcal{A}} \bar{r}(x,a)v(x,a) \\
s.t. \quad & \sum_{x \in \mathcal{S}, a \in \mathcal{A}} (\delta_{xy} - P_{xay})v(x,a) = 0, \quad y \in \mathcal{S} \\
& \sum_{a \in \mathcal{A}} \tau(y,a)v(y,a) + \sum_{x \in \mathcal{S}, a \in \mathcal{A}} (\delta_{xy} - P_{xay})t(x,a) = \beta_y, \quad y \in \mathcal{S} \\
& v(x,a) \geq 0, \; t(x,a) \geq 0 \quad x \in \mathcal{S}, \; a \in \mathcal{A}
\end{aligned}
$$

where $\beta_x > 0$ for $x \in \mathcal{S}$ and $\sum_y \beta_y = 1$. The optimum average expected reward for each initial state is obtained from the dual of this LP.

In the unichain case, the average reward remains constant regarless of the initial state, and the LP reduces to,

$$
\begin{aligned}
\phi_1^* = \quad max \quad & \sum_{x \in \mathcal{S}, a \in \mathcal{A}} \bar{r}(x,a)v(x,a) \\
s.t. \quad & \sum_{x \in \mathcal{S}, a \in \mathcal{A}} (\delta_{xy} - P_{xay})v(x,a) = 0, \quad y \in \mathcal{S} \\
& \sum_{x \in \mathcal{S}, a \in \mathcal{A}} \tau(y,a)v(y,a) = 1 \\
& v(x,a) \geq 0, \quad x \in \mathcal{S}, \; a \in \mathcal{A},
\end{aligned}
$$

with the optimum solution denoted as $\boldsymbol{v}^*$. The optimum pure policy $\boldsymbol{g}^*$ is then obtained by assigning $\boldsymbol{g}_x^*$ in such a way that $v^*(x, \boldsymbol{g}_x^*) > 0$. The optimality equations are given for each state $x$ by,

$$
\zeta_x = \max_a \{\bar{r}(x,a) - g\tau(x,a) + \sum_y P_{xay}\zeta_y\}.
$$

The solution to these equations, $\{\zeta^*, g^*\}$ provides the optimum average expected reward, $\phi_1^* = g^*$.

The constrained problem has been investigated for the average reward SMDPs [6, 7, 16]. Beutler and Ross [6, 7] consider the ratio-average reward with a constraint under a condition stronger than

the unichain condition. In [16], Feinberg examines the problem of maximizing both $\phi_1$ and $\phi_2$ subject to a number of constraints. Under the condition that the initial distribution is fixed, he shows that for both criteria, there exist optimal mixed stationary policies when an associated linear program (LP) is feasible. The mixed policies are defined as policies with an initial one-step randomization applied to a set of pure policies, hence they are not stationary. He provides a linear programming algorithm for the unichain SMDP under both criteria. Average expected reward SMDPs with Borel state and action spaces and unbounded rewards are considered by Schäl [36], Sennott [40], and Luque-Vásquez and Hernández-Lerma [29].

# 6  Expected Time-Average Reward and Variability

The expected time-average reward criterion is similar to the average expected reward criterion. Fatou's lemma immediately implies that $\psi(\boldsymbol{u}) \leq \phi_1(\boldsymbol{u})$ holds for all policies. Baykal-Gürsoy and Gürsoy [4] show that for a large class of policies these two rewards are equal and an $\epsilon$-optimal randomized stationary policy can be obtained for the general (communicating, multichain) SMDP, while such a policy may not exist for the average reward problem [4]. Multiple constraints and the more general expected time-average variability criterion are also discussed. They show that an $\epsilon$-optimal stationary policy can be obtained for the general SMDPs. If $h(x, y) = x - \lambda(x-y)^2$, then the optimal policy is a pure policy. Note that in this case maximizing $\nu(\boldsymbol{u})$ corresponds to maximizing the expected average reward penalized by the expected average variability. A decomposition algorithm to locate the $\epsilon$-optimal stationary policy for both problems is given in [4]. This algorithm utilizes an LP of the form:

$$
\begin{aligned}
max \quad & \sum_{x \in \mathcal{S}, a \in \mathcal{A}} h[\bar{r}(x, a), \sum_{y \in \mathcal{S}, b \in \mathcal{A}} \bar{r}(y, b) v(y, b)] v(x, a) \\
s.t. \quad & \sum_{x \in \mathcal{S}, a \in \mathcal{A}} (\delta_{xy} - P_{xay}) v(x, a) = 0, \quad y \in \mathcal{S} \\
& \sum_{x \in \mathcal{S}, a \in \mathcal{A}} \tau(x, a) v(x, a) = 1 \\
& \sum_{x \in \mathcal{S}, a \in \mathcal{A}} \bar{c}(x, a) v(x, a) \leq \gamma \\
& v(x, a) \geq 0, \quad x \in \mathcal{S}, \ a \in \mathcal{A}.
\end{aligned}
$$

# References

[1] M. Baykal-Gűrsoy. Semi-Markov decision processes: Nonstandard criteria. In *IEEE Regional NY/NJ Conference*, 1994.

[2] M. Baykal-Gűrsoy, T. Altiok, and H. Danhong. Control policies for two-stage production/inventory systems with constrained average cost criterion. *International Journal of Production Research*, 32:2005–2014, 1994.

[3] M. Baykal-Gűrsoy, T. Altiok, and H. Danhong. Look-back policies for two-stage, pull-type production/inventory systems. *Annals of Operations Research*, 48:381–400, 1994. Special Issue on Queueing Networks.

[4] M. Baykal-Gűrsoy and K. Gűrsoy. Semi-Markov decision processes: Nonstandard criteria. *Prob. Eng. Info. Sci.*, 21:635–657, 2007.

[5] M. Baykal-Gűrsoy and K. W. Ross. Variability sensitive Markov decision processes. *Math. Oper. Res.*, 17:558–571, 1992.

[6] F. J. Beutler and K. W. Ross. Time-average optimal constrained semi-Markov decision processes. *Adv. App. Prob.*, 18:341–359, 1986.

[7] F. J. Beutler and K. W. Ross. Uniformization for semi-markov decision processes under stationary policies. *Adv. App. Prob.*, 24:644–656, 1987.

[8] J.S. De Cani. A dynamic programming algorithm for embedded Markov chains when the planning horizon is at infinity. *Mgt. Sci.*, 10(4):716–733, 1964.

[9] E. Çınlar. *Introduction to Stochastic Processes.* Prentice Hall, New Jersey, 1975.

[10] E. V. Denardo. Markov renewal programs with small interest rate. *Annals of Mathematical Statistics.*, 42:477–496, 1971.

[11] E. V. Denardo. Markov renewal programs with small interest rate. *Annals of Mathematical Statistics.*, 42:477–496, 1971.

[12] E. V. Denardo and B. L. Fox. Multichain Markov renewal programs. *SIAM J. Appl. Math.*, 16:468–487, 1968.

[13] A. Federgruen, A. Hordijk, and H. C. Tijms. Denumerable state semi-Markov decison processes with unbounded costs, average cost criterion. *Stoch. Proc. Appl.*, 9, 1979.

[14] A. Federgruen, P. J. Schweitzer, and H. C. Tijms. Denumerable undiscounted semi-Markov decision processes with unbounded rewards. *Math. Oper. Res.*, 8(2), 1983.

[15] A. Federgruen and H. C. Tijms. The optimality equation in average cost denumerable state semi-Markov decison problems, recurrence conditions and algorithms. *J. App. Prob.*, 15, 1978.

[16] E. A. Feinberg. Constrained semi-Markov decision processes with average rewards. *Math. Meth. Oper. Res.*, pages 257–288, 1994.

[17] E. A. Feinberg. Constrained discounted semi-Markov decision processes. In Z. How, J. A. Filar, and A. Chen, editors, *Markov Processes and Controlled Markov Chains*, pages 233–244. Kluwer, Dordrecht, The Netherlands, 2002.

[18] E. A. Feinberg and O. Kella. Optimality of D-policies for an M/G/1 queue with a removable server. *Queueing Systems*, 42:355–376, 2002.

[19] B. Fox. Markov renewal programming by linear fractional programming. *SIAM J. Appl. Math.*, 16:1418–1432, 1966.

[20] B. Hajek. Optimal control of two interacting service stations. *IEEE Transactions on Automatic Control*, 29:491–499, 1984.

[21] D. P. Heyman and M. J. Sobel. *Stochastic Models in Operations Research-Volume II*. McGraw-Hill Book Company, New York, 1984.

[22] Q. Hu and W. Yue. Optimal replacement of a system according to a semi-Markov decision process in a semi-Markov environment. *Optimization Methods and Software*, 18(2):181–196, 2003.

[23] W. S. Jewell. Markov renewal programming I: Formulation, finite return models. *J. Oper. Res.*, 11:938–948, 1963.

[24] W. S. Jewell. Markov renewal programming II: Inifinite return models, example. *J. Oper. Res.*, 11:949–971, 1963.

[25] L. Jianyong and Z. Xiaobo. On average reward semi-Markov decision processes with a general multichain structure. *Math. Oper. Res.*, 29(2):339–352, 2004.

[26] L. C. M. Kallenberg. *Linear Programming and Finite Markovian Control Problems*, volume 148. Mathematical Centre Tracts, Amsterdam, 1983.

[27] S. A. Lippman. Maximal average reward policies for semi-Markov renewal processes with arbitrary state and action spaces. *Ann. Math. Statist.*, 42:1717–1726, 1971.

[28] S. A. Lippman. On dynamic programming with unbounded rewards. *Mgt. Sci.*, 21(11):1225–1233, 1975.

[29] F. Luque-Vásquez and O. Hernández-Lerma. Semi-Markov control models with average costs. *Applicationes Mathematicae*, 26(3):315–331, 1999.

[30] H. Mine and S. Osaki. *Markovian decision processes*. Elsevier, New York, 1970.

[31] M. Pinedo. *Scheduling: Theory, Algorithms, and Systems*. Prentice Hall, Englewood Cliffs, 1995.

[32] M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley and Sons, New York, 1994.

[33] S. Ross. *Applied Probability Models with Optimization Applications*. Holden-Day, San Francisco, 1971.

[34] S. M. Ross. Average cost semi-Markov processes. *J. App. Prob.*, 7:649–656, 1970.

[35] S. M. Ross. *Introduction to Stochastic Dynamic Programming*. Academic Press, 1983.

[36] M. Schǎl. On the second optimality equation for semi-Markov decision models. *Math. Oper. Res.*, 17(2):470–486, 1992.

[37] P. J. Schweitzer. Iterative solution of the functional equations of undiscounted Markov renewal programming. *Journal of Mathematical Analysis and Applications*, 34:495–501, 1971.

[38] P. J. Schweitzer and A. F. Federgruen. The functional equations of undiscounted Markov renewal programming. *Math. Oper. Res.*, 3:308–321, 1978.

[39] P. J. Schweitzer, M. L. Puterman, and K. W. Kindle. Iterative aggregation-disaggregation procedures for solving discounted semi-Markovian reward processes. *Oper. Res.*, 33:589–605, 1985.

[40] L. I. Sennott. Average cost semi-Markov decision processes and the control of queueing systems. *Prob. Eng. Info. Sci.*, 3:247–272, 1989.

[41] L. I Sennott. *Stochastic Dynamic Programming and the Control of Queueing Systems*. Wiley and Sons, New York, 1999.

[42] H. C. Tijms. *Stochastic Modeling and Analysis: A Computational Approach*. Wiley, Chichester, 1986.

[43] H. C. Tijms. *A First Course in Stochastic Models*. Wiley, Chichester, England, 2003.

[44] S. Stidham Jr.and R. R. Weber. A survey of Markov decision models for control of networks of queues. *Queueing Systems*, 13:291–314, 1993.

[45] A. A. Yushkevich. On semi-Markov controlled models with an average reward criterion. *Theory of Probability and Its Applications*, 26:796–802, 1981.