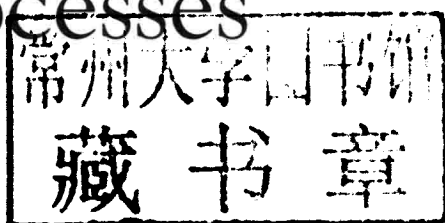


Examples in Markov Decision Processes

A. B. Piunovskiy

Examples in Markov Decision Processes



A. B. Piunovskiy

The University of Liverpool, UK

Published by

Imperial College Press
57 Shelton Street
Covent Garden
London WC2H 9HE

Distributed by

World Scientific Publishing Co. Pte. Ltd.
5 Toh Tuck Link, Singapore 596224
USA office: 27 Warren Street, Suite 401-402, Hackensack, NJ 07601
UK office: 57 Shelton Street, Covent Garden, London WC2H 9HE

British Library Cataloguing-in-Publication Data

A catalogue record for this book is available from the British Library.

Imperial College Press Optimization Series — Vol. 2
EXAMPLES IN MARKOV DECISION PROCESSES

Copyright © 2013 by Imperial College Press

All rights reserved. This book, or parts thereof, may not be reproduced in any form or by any means, electronic or mechanical, including photocopying, recording or any information storage and retrieval system now known or to be invented, without written permission from the Publisher.

For photocopying of material in this volume, please pay a copying fee through the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, USA. In this case permission to photocopy is not required from the publisher.

ISBN 978-1-84816-793-3

Printed in Singapore by Mainland Press Pte Ltd.

Examples in Markov Decision Processes

Imperial College Press Optimization Series

ISSN 2041-1677

Series Editor: Jean Bernard Lasserre (*LAAS-CNRS and Institute of Mathematics, University of Toulouse, France*)

Vol. 1: Moments, Positive Polynomials and Their Applications
by Jean Bernard Lasserre

Vol. 2: Examples in Markov Decision Processes
by A. B. Piunovskiy

Preface

Markov Decision Processes (MDP) is a branch of mathematics based on probability theory, optimal control, and mathematical analysis. Several books with counterexamples/paradoxes in probability [Stoyanov(1997); Szekely(1986)] and in analysis [Gelbaum and Olmsted(1964)] are in existence; it is therefore not surprising that MDP is also replete with unexpected counter-intuitive examples. The main goal of the current book is to collect together such examples. Most of them are based on earlier publications; the remainder are new. This book should be considered as a complement to scientific monographs on MDP [Altman(1999); Bertsekas and Shreve(1978); Hernandez-Lerma and Lasserre(1996a); Hernandez-Lerma and Lasserre(1999); Piunovskiy(1997); Puterman(1994)]. It can also serve as a reference book to which one can turn for answers to curiosities that arise while studying or teaching MDP. All the examples are self-contained and can be read independently of each other. Concerning uncontrolled Markov chains, we mention the illuminating collection of examples in [Suhov and Kelbert(2008)].

A survey of meaningful applications is beyond the scope of the current book. The examples presented either lead to counter-intuitive solutions, or illustrate the importance of conditions in the known theorems. Not all examples are equally simple or complicated. Several examples are aimed at undergraduate students, whilst others will be of interest to professional researchers.

The book has four chapters in line with the four main different types of MDP: the finite-horizon case, infinite horizon with total or discounted loss, and average loss over an infinite time interval. Some basic theoretical statements and proofs of auxiliary assertions are included in the Appendix.

The following notations and conventions will often be used without explanation.

\triangleq means ‘equals by definition’;

\mathbf{C}^∞ is the space of infinitely differentiable functions;

$\mathbf{C}(\mathbf{X})$ is the space of continuous bounded functions on a (topological) space \mathbf{X} ;

$\mathbf{B}(\mathbf{X})$ is the space of bounded measurable functions on a (Borel) space \mathbf{X} ; in discrete (finite or countable) spaces, the discrete topology is usually supposed to be fixed;

$\mathbf{P}(\mathbf{X})$ is the space of probability measures on the (metrizable) space \mathbf{X} , equipped with the weak topology;

If Γ is a subset of space \mathbf{X} then Γ^c is the complement;

$\mathbf{N} = \{1, 2, \dots\}$ is the set of natural numbers; $\mathbf{N}_0 = \mathbf{N} \cup \{0\}$;

\mathbb{R}^N is the N -dimensional Euclidean space; $\mathbb{R} = \mathbb{R}^1$ is the straight line;

$\mathbb{R}^* = [-\infty, +\infty]$ is the extended straight line;

$\mathbb{R}^+ = \{y > 0\}$ is the set of strictly positive real numbers;

$I\{\text{statement}\} = \begin{cases} 1, & \text{if the statement is correct;} \\ 0, & \text{if the statement is false;} \end{cases}$ is the indicator function;

$\delta_a(dy)$ is the Dirac measure concentrated at point a : $\delta_a(\Gamma) = I\{\Gamma \ni a\}$;

If $r \in \mathbb{R}^*$ then $r^+ \triangleq \max\{0, r\}$, $r^- \triangleq \min\{0, r\}$;

$\sum_{i=n}^m f_i \triangleq 0$ and $\prod_{i=n}^m f_i \triangleq 1$ if $m < n$;

$[r]$ is the integer part, the maximal integer i such that $i \leq r$.

Throughout the current book \mathbf{X} is the state space, \mathbf{A} is the action space, $p_t(dy|x, a)$ is the transition probability, $c_t(x, a)$ and $C(x)$ are the loss functions.

Normally, we denote random variables with capital letters (X), small letters (x) being used just for variables, arguments of functions, etc. Bold case (\mathbf{X}) is for spaces. All functions, mappings, and stochastic kernels are assumed to be Borel-measurable unless their properties are explicitly specified.

We say that a function on \mathbb{R}^1 with the values in a Borel space \mathbf{A} is *piece-wise continuous* if there exists a sequence y_i such that $\lim_{i \rightarrow \infty} y_i = \infty$; $\lim_{i \rightarrow -\infty} y_i = -\infty$, this function is continuous on each open interval

(y_i, y_{i+1}) and there exists a right (left) limit as $y \rightarrow y_i + 0$ ($y \rightarrow y_{i+1} - 0$), $i = 0, \pm 1, \pm 2 \dots$. A similar definition is accepted for real-valued piece-wise Lipschitz, continuously differentiable functions.

If \mathbf{X} is a measurable space and ν is a measure on it, then both formulae

$$\int_{\mathbf{X}} f(x) d\nu(x) \quad \text{and} \quad \int_{\mathbf{X}} f(x) \nu(dx)$$

denote the same integral of a real-valued function f with respect to ν .

w.r.t. is the abbreviation for ‘with respect to’, a.s. means ‘almost surely’, and CDF means ‘cumulative distribution function’.

We consider only minimization problems. When formulating theorems and examples published in books (articles) devoted to maximization, we always adjust the statements for our case without any special remarks.

It should be emphasized that the terminology in MDP is not entirely fixed. For example, very often strategies are called policies. There exist several slightly different definitions of a semi-continuous model, and so on.

The author is thankful to Dr.R. Sheen and to Dr.M. Ruck for the proof reading of all the text.

A.B. Piunovskiy

Contents

<i>Preface</i>	v
1. Finite-Horizon Models	1
1.1 Preliminaries	1
1.2 Model Description	3
1.3 Dynamic Programming Approach	5
1.4 Examples	8
1.4.1 Non-transitivity of the correlation	8
1.4.2 The more frequently used control is not better . .	9
1.4.3 Voting	11
1.4.4 The secretary problem	13
1.4.5 Constrained optimization	14
1.4.6 Equivalent Markov selectors in non-atomic MDPs	17
1.4.7 Strongly equivalent Markov selectors in non-atomic MDPs	20
1.4.8 Stock exchange	25
1.4.9 Markov or non-Markov strategy? Randomized or not? When is the Bellman principle violated? . .	27
1.4.10 Uniformly optimal, but not optimal strategy . . .	31
1.4.11 Martingales and the Bellman principle	32
1.4.12 Conventions on expectation and infinities	34
1.4.13 Nowhere-differentiable function $v_t(x)$; discontinuous function $v_t(x)$	38
1.4.14 The non-measurable Bellman function	43
1.4.15 No one strategy is uniformly ε -optimal	44
1.4.16 Semi-continuous model	46

2.	Homogeneous Infinite-Horizon Models: Expected Total Loss	51
2.1	Homogeneous Non-discounted Model	51
2.2	Examples	54
2.2.1	Mixed Strategies	54
2.2.2	Multiple solutions to the optimality equation . . .	56
2.2.3	Finite model: multiple solutions to the optimality equation; conserving but not equalizing strategy .	58
2.2.4	The single conserving strategy is not equalizing and not optimal	58
2.2.5	When strategy iteration is not successful	61
2.2.6	When value iteration is not successful	63
2.2.7	When value iteration is not successful: positive model I	67
2.2.8	When value iteration is not successful: positive model II	69
2.2.9	Value iteration and stability in optimal stopping problems	71
2.2.10	A non-equalizing strategy is uniformly optimal . .	73
2.2.11	A stationary uniformly ε -optimal selector does not exist (positive model)	75
2.2.12	A stationary uniformly ε -optimal selector does not exist (negative model)	77
2.2.13	Finite-action negative model where a stationary uniformly ε -optimal selector does not exist	80
2.2.14	Nearly uniformly optimal selectors in negative models	83
2.2.15	Semi-continuous models and the blackmailer's dilemma	85
2.2.16	Not a semi-continuous model	88
2.2.17	The Bellman function is non-measurable and no one strategy is uniformly ε -optimal	91
2.2.18	A randomized strategy is better than any selector (finite action space)	92
2.2.19	The fluid approximation does not work	94
2.2.20	The fluid approximation: refined model	97
2.2.21	Occupation measures: phantom solutions	101
2.2.22	Occupation measures in transient models	104
2.2.23	Occupation measures and duality	107

2.2.24	Occupation measures: compactness	109
2.2.25	The bold strategy in gambling is not optimal (house limit)	112
2.2.26	The bold strategy in gambling is not optimal (inflation)	115
2.2.27	Search strategy for a moving target	119
2.2.28	The three-way duel (“Truel”)	122
3.	Homogeneous Infinite-Horizon Models: Discounted Loss	127
3.1	Preliminaries	127
3.2	Examples	128
3.2.1	Phantom solutions of the optimality equation	128
3.2.2	When value iteration is not successful: positive model	130
3.2.3	A non-optimal strategy $\hat{\pi}$ for which $v_x^{\hat{\pi}}$ solves the optimality equation	132
3.2.4	The single conserving strategy is not equalizing and not optimal	134
3.2.5	Value iteration and convergence of strategies	135
3.2.6	Value iteration in countable models	137
3.2.7	The Bellman function is non-measurable and no one strategy is uniformly ε -optimal	140
3.2.8	No one selector is uniformly ε -optimal	141
3.2.9	Myopic strategies	141
3.2.10	Stable and unstable controllers for linear systems	143
3.2.11	Incorrect optimal actions in the model with partial information	146
3.2.12	Occupation measures and stationary strategies	149
3.2.13	Constrained optimization and the Bellman principle	152
3.2.14	Constrained optimization and Lagrange multipliers	153
3.2.15	Constrained optimization: multiple solutions	157
3.2.16	Weighted discounted loss and (N, ∞) -stationary selectors	158
3.2.17	Non-constant discounting	160
3.2.18	The nearly optimal strategy is not Blackwell optimal	163
3.2.19	Blackwell optimal strategies and opportunity loss	164

3.2.20	Blackwell optimal and n -discount optimal strategies	165
3.2.21	No Blackwell (Maitra) optimal strategies	168
3.2.22	Optimal strategies as $\beta \rightarrow 1-$ and MDPs with the average loss – I	171
3.2.23	Optimal strategies as $\beta \rightarrow 1-$ and MDPs with the average loss – II	172
4.	Homogeneous Infinite-Horizon Models: Average Loss and Other Criteria	177
4.1	Preliminaries	177
4.2	Examples	179
4.2.1	Why \limsup ?	179
4.2.2	AC-optimal non-canonical strategies	181
4.2.3	Canonical triplets and canonical equations	183
4.2.4	Multiple solutions to the canonical equations in finite models	186
4.2.5	No AC-optimal strategies	187
4.2.6	Canonical equations have no solutions: the finite action space	188
4.2.7	No AC- ε -optimal stationary strategies in a finite state model	191
4.2.8	No AC-optimal strategies in a finite-state semi-continuous model	192
4.2.9	Semi-continuous models and the sufficiency of stationary selectors	194
4.2.10	No AC-optimal stationary strategies in a unichain model with a finite action space	195
4.2.11	No AC- ε -optimal stationary strategies in a finite action model	198
4.2.12	No AC- ε -optimal Markov strategies	199
4.2.13	Singular perturbation of an MDP	201
4.2.14	Blackwell optimal strategies and AC-optimality	203
4.2.15	Strategy iteration in a unichain model	204
4.2.16	Unichain strategy iteration in a finite communicating model	207
4.2.17	Strategy iteration in semi-continuous models	208
4.2.18	When value iteration is not successful	211
4.2.19	The finite-horizon approximation does not work	213

4.2.20	The linear programming approach to finite models	215
4.2.21	Linear programming for infinite models	219
4.2.22	Linear programs and expected frequencies in finite models	223
4.2.23	Constrained optimization	225
4.2.24	AC-optimal, bias optimal, overtaking optimal and opportunity-cost optimal strategies: periodic model	229
4.2.25	AC-optimal and average-overtaking optimal strategies	232
4.2.26	Blackwell optimal, bias optimal, average- overtaking optimal and AC-optimal strategies . .	235
4.2.27	Nearly optimal and average-overtaking optimal strategies	238
4.2.28	Strong-overtaking/average optimal, overtaking optimal, AC-optimal strategies and minimal opportunity loss	239
4.2.29	Strong-overtaking optimal and strong*-overtaking optimal strategies	242
4.2.30	Parrondo's paradox	247
4.2.31	An optimal service strategy in a queueing system	249
Afterword		253
Appendix A Borel Spaces and Other Theoretical Issues		257
A.1	Main Concepts	257
A.2	Probability Measures on Borel Spaces	260
A.3	Semi-continuous Functions and Measurable Selection . . .	263
A.4	Abelian (Tauberian) Theorem	265
Appendix B Proofs of Auxiliary Statements		267
Notation		281
List of the Main Statements		283
<i>Bibliography</i>		285
<i>Index</i>		291

Chapter 1

Finite-Horizon Models

1.1 Preliminaries

A decision maker is faced with the problem of influencing the behaviour of a probabilistic system as it evolves through time. Decisions are made at discrete points in time referred to as *decision epochs* and denoted as $t = 1, 2, \dots, T < \infty$. At each time t , the system occupies a *state* $x \in \mathbf{X}$. The *state space* \mathbf{X} can be either discrete (finite or countably infinite) or continuous (non-empty uncountable Borel subset of a complete, separable metric space, e.g. \mathbb{R}^1). If the state at time t is considered as a random variable, it is denoted by a capital letter X_t ; small letters x_t are just for possible values of X_t . Therefore, the behaviour of the system is described by a *stochastic (controlled) process*

$$X_0, X_1, X_2, \dots, X_T.$$

In case of uncontrolled systems, the theory of Markov processes is well developed: the initial probability distribution for X_0 , $P_0(dx)$, is given, and the dynamics are defined by *transition probabilities* $p_t(dy|x)$. When \mathbf{X} is finite and the process is time-homogeneous, those probabilities form a transition matrix with elements $p(j|i) = P(X_{t+1} = j | X_t = i)$.

In the case of controlled systems, we assume that the *action space* \mathbf{A} is given, which again can be an arbitrary Borel space (including the case of finite or countable \mathbf{A}). As soon as the state X_{t-1} becomes known (equals x_{t-1}), the decision maker must choose an action/control $A_t \in \mathbf{A}$; in general this depends on all the realized values of X_0, X_1, \dots, X_{t-1} along with past actions A_1, A_2, \dots, A_{t-1} . Moreover, that decision can be randomized. The rigorous definition of a control strategy is given in the next section.

As a result of choosing action a at decision epoch t in state x , the decision maker loses $c_t(x, a)$ units, and the system state at the next decision

epoch is determined by the probability distribution $p_t(dy|x, a)$. The function $c_t(x, a)$ is called a *one-step loss*. The *final/terminal* loss equals $C(x)$ when the final state $X_T = x$ is realized.

We assume that the *initial distribution* $P_0(dx)$ for X_0 is given. Suppose a control strategy π is fixed (that is, the rule of choosing actions a_t ; see the next section). Then the random sequence

$$X_0, A_1, X_1, A_2, X_2, \dots, A_T, X_T$$

is well defined: there exists a single probability measure $P_{P_0}^\pi$ on the space of trajectories

$$(x_0, a_1, x_1, a_2, x_2, \dots, a_T, x_T) \in \mathbf{X} \times (\mathbf{A} \times \mathbf{X})^T.$$

For example, if \mathbf{X} is finite and the control strategy is defined by the map $a_t = \varphi_t(x_{t-1})$, then

$$\begin{aligned} P_{P_0}^\pi \{X_0 = i, A_1 = a_1, X_1 = j, A_2 = a_2, X_2 = k, \dots, X_{T-1} = l, A_T = a_T, X_T = m\} \\ = P_0(i)I\{a_1 = \varphi_1(i)\}p_1(j|i, a_1)I\{a_2 = \varphi_2(j)\} \dots p_T(m|l, a_T). \end{aligned}$$

Here and below, $I\{\cdot\}$ is the indicator function; if \mathbf{X} is discrete then transition probabilities $p_t(\cdot|x, a)$ are defined by the values on singletons $p_t(y|x, a)$. The same is true for the initial distribution.

Therefore, for a fixed control strategy π , the *total expected loss* equals $v^\pi = E_{P_0}^\pi[W]$, where

$$W = \sum_{t=1}^T c_t(X_{t-1}, A_t) + C(X_T)$$

is the *total realized loss*. Here and below, $E_{P_0}^\pi$ is the mathematical expectation with respect to probability measure $P_{P_0}^\pi$.

The aim is to find an *optimal* control strategy π^* solving the problem

$$v^\pi = E_{P_0}^\pi \left[\sum_{t=1}^T c_t(X_{t-1}, A_t) + C(X_T) \right] \longrightarrow \inf_{\pi}. \quad (1.1)$$

Sometimes we call v^π the *performance functional*.

Using the dynamic programming approach, under some technical conditions, one can prove the following statement. Suppose function $v_t(x)$ on \mathbf{X} satisfies the following equation

$$\begin{cases} v_T(x) = C(x); \\ v_{t-1}(x) = \inf_{a \in \mathbf{A}} \left\{ c_t(x, a) + \int_{\mathbf{X}} v_t(y) p_t(dy|x, a) \right\} \\ \quad = c_t(x, \varphi_t^*(x)) + \int_{\mathbf{X}} v_t(y) p_t(dy|x, \varphi_t^*(x)); \end{cases} \quad t = T, T-1, \dots, 1. \quad (1.2)$$

Then, the control strategy defined by the map $a_t = \varphi_t^*(x_{t-1})$ solves problem (1.1), i.e. it is optimal; $\inf_{\pi} v^{\pi} = \int_{\mathbf{X}} v_0(x) P_0(dx)$. Therefore, control strategies of the type presented are usually sufficient for solving standard problems. They are called *Markov selectors*.

1.2 Model Description

We now provide more rigorous definitions.

The *Markov Decision Process (MDP)* with a finite horizon is defined by the collection

$$\{\mathbf{X}, \mathbf{A}, T, p, c, C\},$$

where \mathbf{X} and \mathbf{A} are the state and action spaces (Borel); T is the *time horizon*; $p_t(dy|x, a)$, $t = 1, 2, \dots, T$, are measurable stochastic kernels on \mathbf{X} given $\mathbf{X} \times \mathbf{A}$; $c_t(x, a)$ are measurable functions on $\mathbf{X} \times \mathbf{A}$ with values on the extended straight-line $\mathbb{R}^* = [-\infty, +\infty]$; $C(x)$ is a measurable map $C : \mathbf{X} \rightarrow \mathbb{R}^*$. Necessary statements about Borel spaces are presented in Appendix A.

The *space of trajectories* (or *histories*) up to decision epoch t is

$$\mathbf{H}_{t-1} \triangleq \mathbf{X} \times (\mathbf{A} \times \mathbf{X})^{t-1}, \quad t = 1, 2, \dots, T: \quad \mathbf{H} \triangleq \mathbf{X} \times (\mathbf{A} \times \mathbf{X})^T.$$

A *control strategy* $\pi = \{\pi_t\}_{t=1}^T$ is a sequence of measurable stochastic kernels

$$\pi_t(da|x_0, a_1, x_1, \dots, a_{t-1}, x_{t-1}) = \pi_t(da|h_{t-1})$$

on \mathbf{A} , given \mathbf{H}_{t-1} . If a strategy π^m is defined by (measurable) stochastic kernels $\pi_t^m(da|x_{t-1})$ then it will be called a *Markov strategy*. It is called *semi-Markov* if it has the form $\pi_t(da|x_0, x_{t-1})$. A Markov strategy π^{ms} is called *stationary* if none of the kernels $\pi^{\text{ms}}(da|x_{t-1})$ depends on the time t . Very often, stationary strategies are denoted as π^s . If for any $t = 1, 2, \dots, T$ there exists a measurable mapping $\varphi_t(h_{t-1}) : \mathbf{H}_{t-1} \rightarrow \mathbf{A}$ such that $\pi_t(\Gamma|h_{t-1}) = I\{\Gamma \ni \varphi_t(h_{t-1})\}$ for any $\Gamma \in \mathcal{B}(\mathbf{A})$, then the strategy is denoted by the symbol φ and is called a *selector* or *non-randomized strategy*. Selectors of the form $\varphi_t(x_{t-1})$ and $\varphi(x_{t-1})$ are called *Markov* and *stationary* respectively. Stationary semi-Markov strategies and semi-Markov (stationary) selectors are defined in the same way. In what follows, Δ^{All} is the collection of all strategies, Δ^{M} is the set of all Markov strategies, Δ^{MN} is the set of all Markov selectors. In this connection, letter N

corresponds to non-randomized strategies. Further, Δ^S and Δ^{SN} are the sets of all stationary strategies and of all stationary selectors.

We assume that *initial* probability distribution $P_0(dx)$ is fixed. If a control strategy π is fixed too, then there exists a unique probability measure $P_{P_0}^\pi$ on \mathbf{H} such that $P_{P_0}^\pi(\Gamma^X) = P_0(\Gamma^X)$ for $\Gamma \in \mathcal{B}(\mathbf{H}_0) = \mathcal{B}(\mathbf{X})$ and, for all $t = 1, 2, \dots, T$, for $\Gamma^G \in \mathcal{B}(\mathbf{H}_{t-1} \times \mathbf{A})$, $\Gamma^X \in \mathcal{B}(\mathbf{X})$

$$P_{P_0}^\pi(\Gamma^G \times \Gamma^X) = \int_{\Gamma^G} p_t(\Gamma^X | x_{t-1}) P_{P_0}^\pi(dg_t)$$

and

$$P_{P_0}^\pi(\Gamma^H \times \Gamma^A) = \int_{\Gamma^H} \pi_t(\Gamma^A | h_{t-1}) P_{P_0}^\pi(dh_{t-1})$$

for $\Gamma^H \in \mathcal{B}(\mathbf{H}_{t-1})$, $\Gamma^A \in \mathcal{B}(\mathbf{A})$. Here, with some less-than-rigorous notation, we also denote $P_{P_0}^\pi(\cdot)$ the images of $P_{P_0}^\pi$ relative to projections of the types

$$\mathbf{H} \rightarrow \mathbf{H}_{t-1} \times \mathbf{A} \triangleq \mathbf{G}_t, \quad t = 1, 2, \dots, T, \quad \text{and} \quad \mathbf{H} \rightarrow \mathbf{H}_t, \quad t = 0, 1, 2, \dots, T. \quad (1.3)$$

$g_t = (x_0, a_1, x_1, \dots, a_t)$ and $h_t = (x_0, a_1, x_1, \dots, a_t, x_t)$ are the generic elements of \mathbf{G}_t and \mathbf{H}_t . Where they are considered as random elements on \mathbf{H} , we use capital letters G_t and H_t , as usual.

Measures $P_{P_0}^\pi(\cdot)$ on \mathbf{H} are called *strategic* measures; they form space \mathcal{D} .

One can introduce σ -algebras \mathcal{G}_t and \mathcal{F}_t in \mathbf{H} as the pre-images of $\mathcal{B}(\mathbf{G}_t)$ and $\mathcal{B}(\mathbf{H}_t)$ with respect to (1.3). Now the trivial projections

$$(x_0, a_1, x_1, \dots, a_T, x_T) \rightarrow x_t \quad \text{and} \quad (x_0, a_1, x_1, \dots, a_T, x_T) \rightarrow a_t$$

define \mathcal{F} -adapted and \mathcal{G} -adapted stochastic processes $\{X_t\}_{t=0}^T$ and $\{A_t\}_{t=1}^T$ on the *stochastic basis* $(\mathbf{H}, \mathcal{B}(\mathbf{H}), \{\mathcal{F}_0, \mathcal{G}_1, \mathcal{F}_1, \dots, \mathcal{G}_T, \mathcal{F}_T\}, P_{P_0}^\pi)$, which is completed as usual. Note that the process A_t is \mathcal{F} -predictable, and that this property is natural. That is the main reason for considering sequences $(x_0, a_1, x_1, \dots, a_T, x_T)$, not $(x_0, a_0, x_1, \dots, a_{T-1}, x_T)$. The latter notation is also widely used by many authors.

For each $h \in \mathbf{H}$ the (realized) *total loss* equals

$$w(h) = \sum_{t=1}^T c_t(x_{t-1}, a_t) + C(x_T),$$

where we put “ $+\infty$ ” + “ $-\infty$ ” \triangleq “ $+\infty$ ”. The map $W : h \rightarrow w(h)$ defines the random total loss, and the performance of control strategy π is given by $v^\pi = E_{P_0}^\pi[W]$. Here and below,

$$E_{P_0}^\pi[W] \triangleq E_{P_0}^\pi[W^+] + E_{P_0}^\pi[W^-]; \quad \text{“} +\infty \text{”} + \text{“} -\infty \text{”} \triangleq \text{“} +\infty \text{”};$$