

Corrections for
DYNAMIC PROGRAMMING AND
OPTIMAL CONTROL: 2ND and 3RD EDITIONS

by Dimitri P. Bertsekas
Athena Scientific, 2001, 2005

Last Updated: 4/8/08

VOLUME 1 - 3RD EDITION

p. 7 (+19) Change "... after operation B ..." to "... after operation C ..."

p. 213 (+9) This problem is flawed as stated. Replace the statement with the following:

4.27

Consider the quiz contest problem of Example 5.1, where the questions are partitioned in M groups, and there is an order constraint that all the questions in group m must be answered before the questions in group $m + 1$ can be answered. Show that an optimal list can be constructed by ordering the questions within each group in decreasing order of $p_i R_i / (1 - p_i)$. Consider also the problem of optimally ordering the groups in addition to optimally ordering the questions within each group. Show that it is optimal to answer groups in order of decreasing $W / (1 - P)$, where for a given group, W is the expected reward obtained by answering only the questions of that group and in optimal order, and P is the probability of answering all the questions of the group correctly.

p. 270 (-5) Change "... Exercise 5.6 ..." to "... Exercise 5.6 ..."

p. 387 (+20) Change "... all which ..." to "... all of which ..."

p. 428 (-16) Change "... for all i and k ." to "... for all i ."

p. 479 (-14) Change "... Prop. A1 of Appendix A in Vol. II." to "... Section 4.1 of Vol. II."

p. 512 (-3) Change "... with conceptually convenient ..." to "... with the conceptually convenient ..."

VOLUME 2 - 3RD EDITION

p. 198 (-10) Change Prop. 4.1.9 to Prop. 4.2.1

p. 341 (-4) Change $\sum_{i=1}^n$ to $\sum_{j=1}^n$

p. 342 (+7) Change $\sum_{i=1}^n$ to $\sum_{j=1}^n$

p. 351 (-1) Change $(\alpha\lambda)^k$ to $(\alpha\lambda)^t$

p. 369 (-16) Change Eq. (6.66) to read

$$\sum_{t=0}^k \phi(i_t) \tilde{q}(i_{t+1}, r_k) = \sum_{t \leq k, t \in T} \phi(i_t) c(i_{t+1}) + \left(\sum_{t \leq k, t \notin T} \phi(i_t) \phi(i_{t+1})' \right) r_k, \quad (6.66)$$

VOLUME 1 - 2ND EDITION

p. 10 (+21) Change “this true.” to “this is true.”

p. 109 (+4) The expression should read

$$g(x, \mu^*(t, x)) + \nabla_t J^*(t, x) + \nabla_x J^*(t, x)' f(x, \mu^*(t, x)),$$

p. 115 (-8) The equation should read

$$x^*(t) = x(0)e^{at} + \frac{b^2\xi}{2a}(e^{-at} - e^{at}),$$

p. 136, (+10) Change “the optimal $u^*(t)$ ” to “the sine of the slope of the optimal $x^*(t)$ ”

p. 150, (+10) Change “Eq. (4.12)” to “Eq. (4.11)”

p. 156, Fig. 4.2.1 Change $L(y)$ to $G_k(x_k)$

p. 157 (+13) Change $G(x_k)$ to $G_k(x_k)$

p. 160 (-11) Change

$$J_{N-1}(x) = \min \left[cx + G_{N-1}(x), \min_{y>x} [K + cy + G_{N-1}(y)] \right] - cx.$$

to

$$J_{N-1}(x) = \min \left[G_{N-1}(x), \min_{y>x} [K + G_{N-1}(y)] \right] - cx.$$

p. 168 The following is a cleaner version of the three paragraphs starting with the title “Asset Selling”:

Asset Selling

As a first example, consider a person having an asset (say a piece of land) for which he is offered an amount of money from period to period. We assume that the offers, denoted w_0, w_1, \dots, w_{N-1} , are random and independent, and take values within some bounded interval of nonnegative numbers ($w_k = 0$ could correspond to no offer received during the period). If the person accepts an offer, he can invest the money at a fixed rate of interest $r > 0$, and if he rejects the offer, he waits until the next period to consider the next offer. Offers rejected are not renewed, and we assume that the last offer w_{N-1} must be accepted if every prior offer has been rejected. The objective is to find a policy for accepting and rejecting offers that maximizes the revenue of the person at the N th period.

The DP algorithm for this problem can be derived by elementary reasoning. As a modeling exercise, however, we will embed the problem in the framework of the basic problem by specifying the system and cost. We define the state space to be the real line, augmented with an additional state (call it T), which is a *termination state*. By writing that the system is at state $x_k = T$ at some time $k \leq N - 1$, we mean that the asset has already been sold. By writing that the system is at a state $x_k \neq T$ at some time $k \leq N - 1$, we mean that the asset has not been sold as yet and the offer under consideration is equal to x_k (and also equal to the k th offer w_{k-1}). We take $x_0 = 0$ (a fictitious “null” offer). The control space consists of two elements u^1 and u^2 , corresponding to the decisions “sell” and “do not sell,” respectively. We view w_k as the disturbance at time k .

With these conventions, we may write a system equation of the form

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

where the function f_k is defined via the relation

$$x_{k+1} = \begin{cases} T & \text{if } x_k = T, \text{ or if } x_k \neq T \text{ and } u_k = u^1 \text{ (sell),} \\ w_k & \text{otherwise.} \end{cases}$$

Note that a sell decision at time k ($u_k = u^1$) accepts the offer w_{k-1} , and that no explicit sell decision is required to accept the last offer w_{N-1} , as it must be accepted by assumption if the asset has not yet been sold. The corresponding reward function may be written as

$$E_{w_k} \left\{ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k, w_k) \right\}$$

where

$$g_N(x_N) = \begin{cases} x_N & \text{if } x_N \neq T, \\ 0 & \text{otherwise,} \end{cases}$$

$$g_k(x_k, u_k, w_k) = \begin{cases} (1+r)^{N-k} x_k & \text{if } x_k \neq T \text{ and } u_k = u^1 \text{ (sell),} \\ 0 & \text{otherwise.} \end{cases}$$

p. 177 (+15) Change $dp(w)$ to $dP(w)$

p. 187 (+9) and p. 189 (+8) Change $N_k^{-1}C_k$ to N_k^{-1}

p. 217 (+1) Change “If involves” to “It involves”

p. 217 (-7) Change

$$\frac{P(x_1 = \bar{P}, G, G, S)}{P(G, G, S)}$$

to

$$\frac{P(x_1 = \bar{P}, G, G | S)}{P(G, G | S)}$$

p. 218 (+2) Change

$$\frac{P(x_1 = \bar{P}, G, B, S)}{P(G, B, S)}$$

to

$$\frac{P(x_1 = \bar{P}, G, B | S)}{P(G, B | S)}$$

p. 218 (+9) Change

$$\frac{P(x_1 = \bar{P}, G, G, C)}{P(G, G, C)}$$

to

$$\frac{P(x_1 = \bar{P}, G, G | C)}{P(G, G | C)}$$

p. 222 (-9) Change $x'_{N-1}K_{N-1}x_{N-1} | I_{N-2}$ to $E\{x'_{N-1}K_{N-1}x_{N-1} | I_{N-2}, u_{N-2}\}$

p. 244 (+9) The following is a cleaner version of the three pages that start with “**The Conditional State Distribution as a Sufficient Statistic**” title and end just before the “**The Conditional State Distribution Recursion**” title.

The Conditional State Distribution as a Sufficient Statistic

There are many different functions that can serve as sufficient statistics. The identity function $S_k(I_k) = I_k$ is certainly one of them. To obtain another important sufficient statistic, we assume that *the probability distribution of the observation disturbance v_{k+1} depends explicitly only on the immediately preceding state, control, and system disturbance x_k, u_k, w_k , and not on $x_{k-1}, \dots, x_0, u_{k-1}, \dots, u_0, w_{k-1}, \dots, w_0, v_{k-1}, \dots, v_0$* . Under this assumption, it turns out that a sufficient statistic is given by the conditional probability distribution $P_{x_k|I_k}$ of the state x_k , given the information vector I_k . In particular, we will show that for all k and I_k , we have

$$J_k(I_k) = \min_{u_k \in U_k} H_k(P_{x_k|I_k}, u_k) = \bar{J}_k(P_{x_k|I_k}), \quad (5.34)$$

where H_k and \bar{J}_k are appropriate functions.

To this end, we note an important fact that relates to state estimation of discrete-time stochastic systems: the conditional distribution $P_{x_k|I_k}$ can be generated recursively. In particular, it turns out that we can write for all k

$$P_{x_{k+1}|I_{k+1}} = \Phi_k(P_{x_k|I_k}, u_k, z_{k+1}),$$

where Φ_k is some function that can be determined from the data of the problem. Let us postpone a justification of this for the moment, and accept it for the purpose of the following discussion.

We note that to perform the minimization in Eq. (5.32), it is sufficient to know the distribution $P_{x_{N-1}|I_{N-1}}$ together with the distribution $P_{w_{N-1}|x_{N-1},u_{N-1}}$, which is part of the problem data. Thus, the minimization in the right-hand side of Eq. (5.32) is of the form

$$J_{N-1}(I_{N-1}) = \min_{u_{N-1} \in U_{N-1}} H_{N-1}(P_{x_{N-1}|I_{N-1}}, u_{N-1}) = \bar{J}_{N-1}(P_{x_{N-1}|I_{N-1}}),$$

for appropriate functions H_{N-1} and \bar{J}_{N-1} .

We now use induction, i.e., we assume that

$$J_{k+1}(I_{k+1}) = \min_{u_{k+1} \in U_{k+1}} H_{k+1}(P_{x_{k+1}|I_{k+1}}, u_{k+1}) = \bar{J}_{k+1}(P_{x_{k+1}|I_{k+1}}),$$

for appropriate functions H_{k+1} and \bar{J}_{k+1} , and we show that

$$J_k(I_k) = \min_{u_k \in U_k} H_k(P_{x_k|I_k}, u_k) = \bar{J}_k(P_{x_k|I_k}),$$

or appropriate functions H_k and \bar{J}_k .

Indeed, for a given I_k , the expression

$$\min_{u_k \in U_k} E_{x_k, w_k, z_{k+1}} \{g_k(x_k, u_k, w_k) + J_{k+1}(I_{k+1}) \mid I_k, u_k\}$$

in the DP equation (5.33) is written as

$$\min_{u_k \in U_k} E_{x_k, w_k, z_{k+1}} \{g_k(x_k, u_k, w_k) + J_{k+1}(\Phi_k(P_{x_k|I_k}, u_k, z_{k+1})) \mid I_k, u_k\}.$$

In order to calculate the expression being minimized over u_k above, aside from $P_{x_k|I_k}$, we need the joint distribution

$$P(x_k, w_k, z_{k+1} \mid I_k, u_k)$$

or, equivalently,

$$P(x_k, w_k, h_{k+1}(f_k(x_k, u_k, w_k), u_k, v_{k+1}) \mid I_k, u_k).$$

By using Bayes' rule, this distribution can be expressed in terms of $P_{x_k|I_k}$, the given distributions

$$P(w_k \mid x_k, u_k), \quad P(v_{k+1} \mid f_k(x_k, u_k, w_k), u_k, w_k),$$

and the system equation $x_{k+1} = f_k(x_k, u_k, w_k)$. Therefore the expression minimized over u_k can be written as a function of $P_{x_k|I_k}$ and u_k , and the DP equation (5.33) can be written as

$$J_k(I_k) = \min_{u_k \in U_k} H_k(P_{x_k|I_k}, u_k)$$

for a suitable function H_k . Thus the induction is complete and it follows that the distribution $P_{x_k|I_k}$ is a sufficient statistic.

Note that if the conditional distribution $P_{x_k|I_k}$ is uniquely determined by another expression $S_k(I_k)$, that is, $P_{x_k|I_k} = G_k(S_k(I_k))$ for an appropriate function G_k , then $S_k(I_k)$ is also a sufficient statistic. Thus, for example, if we can show that $P_{x_k|I_k}$ is a Gaussian distribution, then the mean and the covariance matrix corresponding to $P_{x_k|I_k}$ form a sufficient statistic.

Regardless of its computational value, the representation of the optimal policy as a sequence of functions of the conditional probability distribution $P_{x_k|I_k}$,

$$\mu_k(I_k) = \bar{\mu}_k(P_{x_k|I_k}), \quad k = 0, 1, \dots, N - 1,$$

is conceptually very useful. It provides a decomposition of the optimal controller in two parts:

- (a) An *estimator*, which uses at time k the measurement z_k and the control u_{k-1} to generate the probability distribution $P_{x_k|I_k}$.
- (b) An *actuator*, which generates a control input to the system as a function of the probability distribution $P_{x_k|I_k}$ (Fig. 5.4.1).

This interpretation has formed the basis for various suboptimal control schemes that separate the controller a priori into an estimator and an actuator and attempt to design each part in a manner that seems “reasonable.” Schemes of this type will be discussed in Chapter 6.

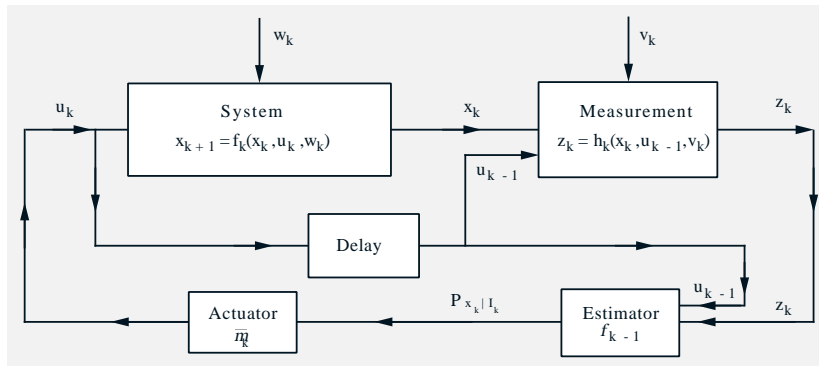


Figure 5.4.1 Conceptual separation of the optimal controller into an estimator and an actuator.

Alternative Perfect State Information Reduction

By using the sufficient statistic $P_{x_k|I_k}$ we can write the DP algorithm in

an alternative form. Using Eq. (5.34), we have for $k < N - 1$

$$\begin{aligned} \bar{J}_k(P_{x_k|I_k}) = \min_{u_k \in U_k} & \left[E_{x_k, w_k, z_{k+1}} \{g_k(x_k, u_k, w_k) \right. \\ & \left. + \bar{J}_{k+1}(\Phi_k(P_{x_k|I_k}, u_k, z_{k+1})) \mid I_k, u_k \right]. \end{aligned} \quad (5.35)$$

In the case where $k = N - 1$, we have

$$\begin{aligned} \bar{J}_{N-1}(P_{x_{N-1}|I_{N-1}}) \\ = \min_{u_{N-1} \in U_{N-1}} & \left[E_{x_{N-1}, w_{N-1}} \left\{ g_N(f_{N-1}(x_{N-1}, u_{N-1}, w_{N-1})) \right. \right. \\ & \left. \left. + g_{N-1}(x_{N-1}, u_{N-1}, w_{N-1}) \mid I_{N-1}, u_{N-1} \right\} \right]. \end{aligned} \quad (5.36)$$

This DP algorithm yields the optimal cost as

$$J^* = E_{z_0} \{ \bar{J}_0(P_{x_0|z_0}) \},$$

where \bar{J}_0 is obtained by the last step, and the probability distribution of z_0 is obtained from the measurement equation $z_0 = h_0(x_0, v_0)$ and the statistics of x_0 and v_0 .

By observing the form of Eq. (5.35), we note that it has the standard DP structure, except that $P_{x_k|I_k}$ plays the role of the “state.” Indeed the role of the “system” is played by the recursive estimator of $P_{x_k|I_k}$,

$$P_{x_{k+1}|I_{k+1}} = \Phi_k(P_{x_k|I_k}, u_k, z_{k+1}),$$

and this system fits the framework of the basic problem (the role of control is played by u_k and the role of the disturbance is played by z_{k+1}). Furthermore, the controller can calculate (at least in principle) the state $P_{x_k|I_k}$ of this system at time k , so perfect state information prevails. Thus the alternate DP algorithm (5.34)-(5.35) may be viewed as the DP algorithm of the perfect state information problem that involves the above system, whose state is $P_{x_k|I_k}$, and an appropriately reformulated cost function. In the absence of perfect knowledge of the state, *the controller can be viewed as controlling the “probabilistic state” $P_{x_k|I_k}$ so as to minimize the expected cost-to-go conditioned on the information I_k available.*

p. 254 (-1) Change a_k to α_k

p. 255 (+6) Change $1 - \frac{1}{C}$ to $1 - \frac{I}{C}$

p. 261 (+8) Change “Exercise 5.6” to “Exercise 5.3”

p. 262 (+15) Change $y'_0 K_k y_0$ to $y'_0 K_0 y_0$

p. 266 (-5) Delete part (e) (it is correct only for open-loop policies)

- p. 316 (+16)** Change “of question” to “of questions”
- p. 349 (+15)** Change “Section 1.3” to “Section 2.3”
- p. 353 (+12)** Change “CEC is 1.” to “CEC with nominal values $\bar{w}_0 = \bar{w}_1 = 0$ is 1.”
- p. 357 (+13)** Replace Exercise 6.13 by the following:

6.13 (Discretization of Convex Problems)

Consider a problem with state space S , for all k , where S is a convex subset of \mathfrak{R}^n . Suppose that $\hat{S} = \{y_1, \dots, y_M\}$ is a finite subset of S such that S is the convex hull of \hat{S} , and consider a one-step lookahead policy based on approximated cost-to-go functions $\tilde{J}_0, \tilde{J}_1, \dots, \tilde{J}_N$ defined as follows:

$$\tilde{J}_N(x) = g_N(x), \quad \forall x \in S,$$

and for $k = 1, \dots, N - 1$,

$$\tilde{J}_k(x) = \min \left\{ \sum_{i=1}^M \lambda_i \hat{J}_k(y_i) \mid \sum_{i=1}^M \lambda_i y_i = x, \sum_{i=1}^M \lambda_i = 1, \lambda_i \geq 0, i = 1, \dots, M \right\},$$

where $\hat{J}_k(x)$ is defined by

$$\hat{J}_k(x) = \min_{u \in U_k(x)} E \left\{ g_k(x, u, w_k) + \tilde{J}_{k+1}(f_k(x, u, w_k)) \right\}, \quad \forall i = 1, \dots, M.$$

Thus \tilde{J}_k is obtained from \tilde{J}_{k+1} as a “grid-based” convex piecewise linear approximation to \hat{J}_k based on the M values

$$\hat{J}_k(y_1), \dots, \hat{J}_k(y_M).$$

Assume that the cost functions g_k and the system functions f_k are such that the function \hat{J}_k is real-valued and convex over S whenever \tilde{J}_{k+1} is real-valued and convex over S . Show that the cost-to-go functions $\bar{J}_k(x_k)$ corresponding to the one-step lookahead policy satisfies for all $x \in S$

$$\bar{J}_k(x) \leq \hat{J}_k(x) \leq \tilde{J}_k(x), \quad k = 0, 1, \dots, N - 1.$$

Hint: Use Prop. 6.3.1.

- p. 373 (-1)** Replace $N_*(j)$ by $N^*(j)$
- p. 374 (+3), (+8)** Replace $N_*(j)$ by $N^*(j)$
- p. 396 (+6)** The expression should read

$$\bar{\tau}_i(u) = \sum_{j=1}^n \int_0^\infty \tau dQ_{ij}(\tau, u),$$

p. 402 (+15), (+19), (-5) Change $g(i, u)$ to $G(i, u)$

p. 402 (-5) After Eq. (7.54), add the following sentence: If there is an “instantaneous” one-stage cost $\hat{g}(i, u)$, the term $G(i, u)$ should be replaced by $\hat{g}(i, u) + G(i, u)$ in this equation.

p. 408 (+2) Change (7.6) to (7.17)

p. 451 (+12) Change $C_k N_k^{-1}$ to $C'_k N_k^{-1}$

p. 476 (-11) Change

$$P_{d_1} \preceq P_{d_2} \text{ if and only if } E\{U(f(d_1, n)) \mid d_1\} \leq \{U(f(d_2, n)) \mid d_2\}.$$

to

$$P_{d_1} \preceq P_{d_2} \text{ if and only if } E\{U(f(d_1, n)) \mid d_1\} \leq E\{U(f(d_2, n)) \mid d_2\}.$$

VOLUME 2 - 2ND EDITION

p. 7 (+19) Change “operation D can be performed only after operation B has been performed” to “operation D can be performed only after operation C has been performed”

p. 50 (+17) Change “Tsitsiklis” to “Castanon”

p. 64 (+13) Change the first four lines of the proof as follows:

Proof: In view of Eqs. (1.59) and (1.65), existence of a PPR policy is equivalent to having, for all i ,

$$\max \left\{ M, \max_{j \neq i} L^j(x, M, J) \right\} \geq L^i(x, M, J), \quad \text{for all } x \text{ with } x^i \in S^i, \quad (1.66)$$

$$M \leq L^i(x, M, J), \quad \text{for all } x \text{ with } x^i \notin S^i, \quad (1.67)$$

p. 70 (+10) Change p_i to P_i

p. 74 (+1) Change “... , [VeP84], and Verd’u and Poor [VeP87].” to “... , and Verd’u and Poor [VeP84], [VeP87].”

p. 80 (+17,+19,27) Change 1 to s

p. 84 (+10) Change the equation to

$$T_{\bar{\mu}}J = TJ. \quad (1.87)$$

p. 94 (-1) Change equation to

$$J_{\mu}(1) = -(1 - u^2)u + (1 - u^2)J_{\mu}(1)$$

p. 95 (+2) Change equation to

$$J_{\mu}(1) = -\frac{1 - u^2}{u}.$$

p. 123 (+6) Change “thrtwe” to (2.33)

p. 136 (+3) Change “an optimal proper policy” to “a proper policy”

p. 136 (+11) Change “optimal and proper.” to “optimal.”

p. 136 (-10) Change the last four lines of the hint to:

By taking limit superior as $N \rightarrow \infty$, we obtain $J^* \geq TJ^* \geq J_{\mu}$. Therefore, μ is an optimal policy and we have $J^* = TJ^*$. For the rest of the proof follow the line of proof of Prop. 2.1.2.

- p. 181 (-2)** Change “discount factor” to “discount factor with $\alpha < 2$ ”
- p. 190 (-1)** Change “no optimal policy (stationary or not).” to “no optimal stationary policy.”
- p. 206 (-1)** Change “ N_1 ” to “ $N - 1$ ”
- p. 208 (-12)** Change Exercise number to 4.32.
- p. 234 (+20)** Change “Prop. 4.3.3” to “Prop. 4.2.6”
- p. 235 (+14)** Change “the probabilities $q(i, u)$, $u \in U(i)$.” to “some probabilities.”
- p. 247 (-4)** Change “for some $\beta > 0$ ” to “for some $\beta \in (0, 1)$ ”
- p. 264 (-6)** Change a_k to α_k
- p. 265 (+3)** Change $1 - \frac{1}{C}$ to $1 - \frac{I}{C}$