

**e - c o m p a n i o n**

ONLY AVAILABLE IN ELECTRONIC FORM

Electronic Companion—“Information Relaxations and Duality in Stochastic Dynamic Programs” by David B. Brown, James E. Smith, and Peng Sun, *Operations Research*, DOI 10.1287/opre.1090.0796.

---

# Online Appendices for Information Relaxations and Duality in Stochastic Dynamic Programs

David B. Brown, James E. Smith, and Peng Sun  
Fuqua School of Business  
Duke University

## A. Proofs and Supplemental Discussion

### A.1. Proof of Theorem 2.1

*Proof.* By weak duality, the right side of (7) is greater than or equal to the left side. To establish strong duality, we need to show that if the left side is bounded, there exists a  $z^*$  that obtains equality. Let  $z^*(a) = r(a) - v^*$  where  $v^*$  is the optimal value of the primal DP (1). To see that this  $z$  is dual feasible, note that for  $\alpha_F \in \mathcal{A}_F$ ,  $\mathbb{E}[z^*(\alpha_F)] = \mathbb{E}[r(\alpha_F)] - v^*$ . Because  $\mathbb{E}[r(\alpha_F)] \leq v^*$  (by definition of  $v^*$  as the supremum over policies in  $\mathcal{A}_F$ ), we know that  $\mathbb{E}[z^*(\alpha)] \leq 0$  for all  $\alpha_F \in \mathcal{A}_F$ ; thus  $z^*$  is dual feasible. With this penalty, the penalized objective  $r(a) - z(a)$  is equal to  $v^*$  for all  $a$  and hence, for any policy  $\alpha$  (including those in  $\mathcal{A}_G$ ),  $\mathbb{E}[r(\alpha) - z^*(\alpha)] = v^*$ . This implies that  $z^*$  achieves equality in (7).

If the primal problem is unbounded, by weak duality, the dual problem must also be unbounded.  $\square$

### A.2. Proof of Theorem 2.2

*Proof.* We first consider sufficiency. Consider any  $\alpha_F^* \in \mathcal{A}_F$  and  $z^* \in \mathcal{Z}_F$  and suppose (8) holds and  $\mathbb{E}[z^*(\alpha_F^*)] = 0$ . Then we can rewrite the dual problem with this penalty as

$$\begin{aligned} \sup_{\alpha_G \in \mathcal{A}_G} \mathbb{E}[r(\alpha) - z^*(\alpha)] &= \mathbb{E}[r(\alpha_F^*) - z^*(\alpha_F^*)] && \text{(using (8))} \\ &= \mathbb{E}[r(\alpha_F^*)] && \text{(since } \mathbb{E}[z^*(\alpha_F^*)] = 0\text{)}. \end{aligned}$$

Then, by weak duality,  $\alpha_F^*$  and  $z^*$  must be optimal.

To show necessity, first note that for any  $\alpha_F^* \in \mathcal{A}_F$  and  $z^* \in \mathcal{Z}_F$ , we have:

$$\begin{aligned} \sup_{\alpha_G \in \mathcal{A}_G} \mathbb{E}[r(\alpha_G) - z^*(\alpha_G)] &\geq \sup_{\alpha_F \in \mathcal{A}_F} \mathbb{E}[r(\alpha_F) - z^*(\alpha_F)] && \text{(because } \mathcal{A}_F \subseteq \mathcal{A}_G\text{)} \\ &\geq \mathbb{E}[r(\alpha_F^*) - z^*(\alpha_F^*)] && \text{(because } \alpha_F^* \in \mathcal{A}_F\text{)} \\ &\geq \mathbb{E}[r(\alpha_F^*)] && \text{(because } z^* \in \mathcal{Z}_F\text{)}. \end{aligned}$$

If  $\alpha_F^* \in \mathcal{A}_F$  and  $z^* \in \mathcal{Z}_F$  are primal and dual optimal (respectively), then by the strong duality theorem, the first and last terms above are equal, so the intervening inequalities must hold with equality and we have  $\mathbb{E}[z^*(\alpha_F^*)] = 0$  and (8).  $\square$

### A.3. Proof of Proposition 2.1

*Proof.* The first inequality in (9) follows from applying the weak duality result (Lemma 2.1) with the restricted policy space  $\mathcal{S}$  in place of the full policy space  $\mathcal{A}$ . Note that, by definition, any dual feasible penalty  $z$  for the original problem with  $\mathcal{A}$  satisfies  $\mathbb{E}[z(\alpha_F)] \leq 0$  for all  $\alpha_F$  in  $\mathcal{A}_F$ . Since  $\mathcal{S} \subseteq \mathcal{A}$ , any such penalty will also be dual feasible with a restricted policy space, i.e.,  $\mathbb{E}[z(\alpha_F)] \leq 0$  for all  $\alpha_F$  in  $\mathcal{S}_F$ . This set of dual feasible penalties in the restricted policy space is larger than the set of dual feasible penalties in original space. Thus this first inequality holds on the larger set of penalties that are dual feasible with the restricted penalties.

The second inequality in (9) follows from the fact that  $\mathcal{S}_G \subseteq \mathcal{A}_G$ .  $\square$

### A.4. Proof of Proposition 2.2

Before proving Proposition 2.2, we first establish a lemma that is used in the proof of this proposition as well as some later results.

**Lemma A.1.** *Let  $z_t^G(a) = \mathbb{E}[w_t(a)|\mathcal{G}_t]$ . If  $w_t(a)$  depends on the first  $t+1$  actions in  $a$  and  $\alpha_G$  is  $\mathbb{G}$ -adapted, then  $z_t^G(\alpha_G) = \mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t]$ , almost everywhere.*

Note that this result need not hold for policies that are not  $\mathbb{G}$ -adapted: In  $z_t^{\mathbb{G}}(\alpha_G)$ , the policy  $\alpha_G$  selects actions  $a$  for a given  $\omega$  and  $z_t^{\mathbb{G}}(\alpha_G)$  takes on the corresponding value of  $\mathbb{E}[w_t(a)|\mathcal{G}_t]$ . That is, we calculate the “ $\mathbb{G}$ -average” in the conditional expectation first and then select averaged values. In  $\mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t]$  we select values  $w_t(a)$  according to the policy  $\alpha_G$  first and then calculate the averages in the conditional expectations. In these terms, the lemma says that if  $\alpha_G$  is  $\mathbb{G}$ -adapted,  $\mathbb{G}$ -averaging then selecting is equivalent to selecting then  $\mathbb{G}$ -averaging.

*Proof.* Consider a  $\mathbb{G}$ -adapted policy  $\alpha_G$  and pick an  $\omega^0$  in  $\Omega$  and let  $a^0 = \alpha_G(\omega^0)$ . Define  $H^0$  as the set of  $\omega$  such that the first  $t+1$  actions in  $a$  match those of  $a^0$ ; note that  $\omega^0$  is in  $H^0$ . Because  $\alpha_G$  is  $\mathbb{G}$ -adapted, we know  $H^0 \in \mathcal{G}_t$ . For any set  $H \in \mathcal{G}_t$  such that  $H \subseteq H^0$ , we have

$$\int_H z_t^{\mathbb{G}}(\alpha_G) d\mathbb{P} = \int_H \mathbb{E}[w_t(a^0)|\mathcal{G}_t] d\mathbb{P} = \int_H w_t(a^0) d\mathbb{P}.$$

The first equality follows from the definition of  $z_t^{\mathbb{G}}(\alpha_G)$  taking into account the fact  $H \subseteq H^0$  and that, using  $\alpha_G$ , all  $\omega$  in  $H^0$  select the same actions  $(a_0, \dots, a_t)$  as  $a^0$ . The second equality follows from the definition of conditional expectations (see, e.g., Billingsley 1986, p. 466)<sup>1</sup> using the fact that  $H \in \mathcal{G}_t$ . Similarly, for  $H \in \mathcal{G}_t$  such that  $H \subseteq H^0$ , we have

$$\int_H \mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t] d\mathbb{P} = \int_H w_t(\alpha_G) d\mathbb{P} = \int_H w_t(a^0) d\mathbb{P}.$$

Here we first use the definition of the conditional expectations and then use the fact that, under  $\alpha_G$ , all  $\omega$  in  $H^0$  select the same actions  $(a_0, \dots, a_t)$  as  $a^0$ . Thus for  $H \in \mathcal{G}_t$  such that  $H \subseteq H^0$ , we have

$$\int_H z_t^{\mathbb{G}}(\alpha_G) d\mathbb{P} = \int_H \mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t] d\mathbb{P}. \quad (25)$$

We now show (25) is sufficient to ensure that  $z_t^{\mathbb{G}}(\alpha_G) = \mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t]$  for almost all  $\omega \in H$ . Note that  $z_t^{\mathbb{G}}(\alpha_G)$  and  $\mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t]$  are both  $\mathcal{G}_t$ -measurable; thus  $f = z_t^{\mathbb{G}}(\alpha_G) - \mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t]$  is also  $\mathcal{G}_t$ -measurable. Let  $H^+$  be the subset of  $H^0$  where  $f$  is strictly positive; that is the subset of  $H^0$  where  $z_t^{\mathbb{G}}(\alpha_G) > \mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t]$ . Because  $f$  is  $\mathcal{G}_t$ -measurable,  $H^+ \in \mathcal{G}_t$ . Then from (25), we know that  $\int_{H^+} f d\mathbb{P} = 0$ , which (because  $f > 0$  on  $H^+$ ) implies that  $H^+$  has measure 0 (see, e.g., Billingsley 1986 p. 466). We can similarly define a set  $H^-$  where  $f < 0$  and conclude that  $H^-$  has measure 0. Thus, we can conclude that  $f = 0$  or  $z_t^{\mathbb{G}}(\alpha_G) = \mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t]$  holds almost surely on  $H^0$ . Since we can construct such a set  $H^0$  containing  $\omega^0$  for any  $\omega^0$ , we can conclude that  $z_t^{\mathbb{G}}(\alpha_G) = \mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t]$  for almost all  $\omega$ .  $\square$

We now turn to the proof of Proposition 2.2.

*Proof.* (i) We can write  $z_t(a) = z_t^{\mathbb{G}}(a) - z_t^{\mathbb{F}}(a)$  where  $z_t^{\mathbb{G}}(a) = \mathbb{E}[w_t(a)|\mathcal{G}_t]$  and  $z_t^{\mathbb{F}}(a) = \mathbb{E}[w_t(a)|\mathcal{F}_t]$ . Any  $\alpha_F$  that is  $\mathbb{F}$ -adapted is also  $\mathbb{G}$ -adapted because  $\mathbb{G}$  is a relaxation of  $\mathbb{F}$ . We can then write

$$\mathbb{E}[z_t^{\mathbb{G}}(\alpha_F)|\mathcal{F}_t] = \mathbb{E}[\mathbb{E}[w_t(\alpha_F)|\mathcal{G}_t]|\mathcal{F}_t] = \mathbb{E}[w_t(\alpha_F)|\mathcal{F}_t]. \quad (26)$$

Here the first equality follows from Lemma A.1 and the second from the “law of iterated expectations” (see e.g., Billingsley 1986, p. 470) since  $\mathcal{F}_t \subseteq \mathcal{G}_t$ . Then, using (26) and Lemma A.1 again, we have

$$\mathbb{E}[z_t(\alpha_F)|\mathcal{F}_t] = \mathbb{E}[z_t^{\mathbb{G}}(\alpha_F)|\mathcal{F}_t] - \mathbb{E}[z_t^{\mathbb{F}}(\alpha_F)|\mathcal{F}_t] = \mathbb{E}[w_t(\alpha_F)|\mathcal{F}_t] - \mathbb{E}[w_t(\alpha_F)|\mathcal{F}_t] = 0. \quad (27)$$

This establishes the first part of claim (i). Applying the law of iterated expectations then implies  $\mathbb{E}[z_t(\alpha_F)] = \mathbb{E}[\mathbb{E}[z_t(\alpha_F)|\mathcal{F}_t]] = 0$ . Summing these these over time implies that  $\mathbb{E}[z(\alpha_F)] = 0$ , as stated in the second part of claim (i).

(ii) The fact that  $z$  is adapted to  $\mathbb{G}$  follows from the definition of conditional expectations: for any  $a$ ,  $\mathbb{E}[w_t(a)|\mathcal{G}_t]$  is  $\mathcal{G}_t$ -measurable and  $\mathbb{E}[w_t(a)|\mathcal{F}_t]$  is  $\mathcal{F}_t$ -measurable and hence  $\mathcal{G}_t$ -measurable, since  $\mathbb{G}$  is a

<sup>1</sup>Billingsley, Patrick (1986). Probability and Measure. John Wiley and Sons, New York.

relaxation of  $\mathbb{F}$ . Then  $z_t(a) = \mathbb{E}[w_t(a)|\mathcal{G}_t] - \mathbb{E}[w_t(a)|\mathcal{F}_t]$  is  $\mathcal{G}_t$ -measurable. The fact that  $z_t(a)$  depends only on the first  $t+1$  actions  $(a_0, \dots, a_t)$  of  $a$  follows from  $w_t(a)$  having this same property.  $\square$

### A.5. Proof of Theorem 2.3

*Proof.* The fact that  $z^*$  is dual feasible follows from Proposition 2.2 and the fact that it is optimal for the dual problem follows from the inductive argument in the text preceding the statement of the theorem. The fact that any  $\alpha_F^* \in \mathcal{A}_{\mathbb{F}}$  that is optimal for the primal problem is also optimal for the dual problem then follows by the complementary slackness result, Theorem 2.2.

To establish the last part of the theorem, let us abuse notation a bit and write  $V_t(a)$  in place of  $V_t(a_0, \dots, a_t)$  with the understanding that subsequence of actions  $(a_0, \dots, a_t)$  is selected from the full sequence of actions  $a$ . If  $\alpha_F^* \in \mathcal{A}_{\mathbb{F}}$  is optimal for the primal problem, we then have

$$\begin{aligned} r(\alpha_F^*) - z^*(\alpha_F^*) &= \sum_{t=0}^T r_t(\alpha_F^*) - z_t^*(\alpha_F^*) \\ &= \sum_{t=0}^T r_t(\alpha_F^*) + \mathbb{E}[V_{t+1}(\alpha_F^*)|\mathcal{F}_t] - \mathbb{E}[V_{t+1}(\alpha_F^*)|\mathcal{G}_t] \\ &= \sum_{t=0}^T V_t(\alpha_F^*) - \mathbb{E}[V_{t+1}(\alpha_F^*)|\mathcal{G}_t] \text{ (almost everywhere)}. \end{aligned} \quad (28)$$

The first equality is simply the definition of  $r$  and  $z^*$  and the second equality follows from Lemma A.1 (using the fact that  $\alpha_F^* \in \mathcal{A}_{\mathbb{F}}$ ). The final equality follows from the definition of  $V_t$ , i.e., we have  $V_t(\alpha_F^*) = r_t(\alpha_F^*) + \mathbb{E}[V_{t+1}(\alpha_F^*)|\mathcal{F}_t]$  (almost everywhere); note this equality may fail on a set of measure zero for an optimal policy  $\alpha_F^*$ .

If  $\mathbb{G}$  is the perfect information relaxation, then  $\mathbb{E}[V_{t+1}(\alpha_F^*)|\mathcal{G}_t] = V_{t+1}(\alpha_F^*)$ , adjacent terms in (28) cancel and (28) reduces to  $V_0(\alpha_F^*) - V_{T+1}(\alpha_F^*)$ . Here  $V_{T+1}$  was defined to be 0 and  $V_0(\alpha_F^*)$  is equal to  $\mathbb{E}[r(\alpha_F^*)]$ . Thus, with a perfect information relaxation, we have  $r(\alpha_F^*) - z^*(\alpha_F^*) = \mathbb{E}[r(\alpha_F^*)]$ , almost everywhere.  $\square$

### A.6. Proof of Proposition 2.3

*Proof.* (i) Let  $z_t^{\mathbb{F}}(a) = \mathbb{E}[w_t(a)|\mathcal{F}_t]$ . The inequality in (12) can be established as follows

$$\begin{aligned} \sup_{\alpha_G \in \mathcal{A}_{\mathbb{G}^1}} \mathbb{E}[r(\alpha_G) - z^1(\alpha_G)] &= \sup_{\alpha_G \in \mathcal{A}_{\mathbb{G}^1}} \mathbb{E} \left[ r(\alpha_G) - \sum_t (\mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t^1] - z_t^{\mathbb{F}}(\alpha_G)) \right] \\ &\leq \sup_{\alpha_G \in \mathcal{A}_{\mathbb{G}^2}} \mathbb{E} \left[ r(\alpha_G) - \sum_t (\mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t^1] - z_t^{\mathbb{F}}(\alpha_G)) \right] \\ &= \sup_{\alpha_G \in \mathcal{A}_{\mathbb{G}^2}} \mathbb{E} \left[ r(\alpha_G) - \sum_t (\mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t^2] - z_t^{\mathbb{F}}(\alpha_G)) \right] \\ &= \sup_{\alpha_G \in \mathcal{A}_{\mathbb{G}^2}} \mathbb{E}[r(\alpha_G) - z^2(\alpha_G)] \end{aligned}$$

The first equality follows from the definition of the penalty  $z^1$  and Lemma A.1. The inequality follows from the fact that  $\mathcal{A}_{\mathbb{G}^1} \subseteq \mathcal{A}_{\mathbb{G}^2}$  when  $\mathbb{G}^2$  is a relaxation of  $\mathbb{G}^1$ . The next equality follows from the fact that  $\mathbb{E}[\mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t^2]] = \mathbb{E}[\mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t^1]]$  which we will establish shortly. The final equality follows from the definition of the penalty  $z^2$  and Lemma A.1. To see that  $\mathbb{E}[\mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t^2]] = \mathbb{E}[\mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t^1]]$ , note that  $\mathbb{G}^2$  being a relaxation of  $\mathbb{G}^1$  implies that  $\mathcal{G}_t^1 \subseteq \mathcal{G}_t^2$ . Then, by the law of iterated expectations, we have  $\mathbb{E}[\mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t^2]] = \mathbb{E}[\mathbb{E}[\mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t^2]|\mathcal{G}_t^1]] = \mathbb{E}[\mathbb{E}[w_t(\alpha_G)|\mathcal{G}_t^1]]$ .

(ii) For any two dual feasible penalties  $z^1$  and  $z^2$  and information relaxation  $\mathbb{G}$ , we have

$$\sup_{\alpha_G \in \mathcal{A}_{\mathbb{G}}} \mathbb{E}[r(\alpha_G) - z^1(\alpha_G)] = \sup_{\alpha_G \in \mathcal{A}_{\mathbb{G}}} \mathbb{E}[r(\alpha_G) - z^2(\alpha_G) + (z^2(\alpha_G) - z^1(\alpha_G))] \quad (29)$$

$$\leq \sup_{\alpha_G \in \mathcal{A}_{\mathbb{G}}} \mathbb{E}[r(\alpha_G) - z^2(\alpha_G)] + \sup_{\alpha_G \in \mathcal{A}_{\mathbb{G}}} \mathbb{E}[z^2(\alpha_G) - z^1(\alpha_G)] \quad (30)$$

Rearranging this yields the inequality on the right in (13). Interchanging  $z^1$  and  $z^2$  and multiplying through by  $-1$  yields the inequality on the left in (13).

(iii) The proof here follows the proof of Proposition 2.2, except we now take  $z_t^{\mathbb{F}}(a) = \mathbb{E}[w_t(a)|\mathcal{F}'_t]$  and must show that  $\mathbb{E}[z_t^{\mathbb{F}}(\alpha_F)|\mathcal{F}_t] = \mathbb{E}[w_t(\alpha_F)|\mathcal{F}_t]$  for any  $\mathbb{F}$ -adapted  $\alpha_F$ . This can be established as follows:

$$\mathbb{E}[z_t^{\mathbb{F}}(\alpha_F)|\mathcal{F}_t] = \mathbb{E}[\mathbb{E}[w_t(\alpha_F)|\mathcal{F}'_t]|\mathcal{F}_t] = \mathbb{E}[w_t(\alpha_F)|\mathcal{F}_t]. \quad (31)$$

Here the first equality follows from Lemma A.1 (since  $\alpha_F$  being  $\mathbb{F}$ -adapted implies  $\alpha_F$  is also  $\mathbb{F}'$ -adapted) and the second equality follows from the law of iterated expectations since  $\mathcal{F}_t \subseteq \mathcal{G}_t$ . The rest of the proof then proceeds as in the proof of Proposition 2.2.

(iv) Suppose  $\alpha_G^* \in \mathcal{A}_{\mathbb{G}}$  is an optimal solution for the left side of (14). Then we have

$$\begin{aligned} \sup_{\alpha_G \in \mathcal{A}_{\hat{\mathbb{G}}}} \mathbb{E}[r(\alpha_G) - \hat{z}(\alpha_G)] &\geq \mathbb{E}[r(\alpha_G^*) - \hat{z}(\alpha_G^*)] \\ &= \mathbb{E}\left[r(\alpha_G^*) - \sum_{t=0}^T \hat{z}_t(\alpha_G^*)\right] \\ &= \mathbb{E}\left[r(\alpha_G^*) - \sum_{t=0}^T \mathbb{E}[\hat{z}_t(\alpha_G^*)|\mathcal{G}_t]\right] \\ &= \mathbb{E}\left[r(\alpha_G^*) - \sum_{t=0}^T z_t(a)\right] \\ &= \mathbb{E}[r(\alpha_G^*) - z(\alpha_G^*)] \\ &= \sup_{\alpha_G \in \mathcal{A}_{\mathbb{G}}} \mathbb{E}[r(\alpha_G) - z(\alpha_G)] \end{aligned}$$

The first inequality follows from the fact that  $\alpha_G^* \in \mathcal{A}_{\hat{\mathbb{G}}}$  (since  $\mathbb{G} \subseteq \hat{\mathbb{G}}$ ). The next two equalities follow from the definition of  $\hat{z}$  and the law of iterated expectations, respectively. The third equality follows from the estimate being unbiased and Lemma A.1: the estimate being unbiased means  $z_t(a) = \mathbb{E}[\hat{z}_t(a)|\mathcal{G}_t]$  and, since  $\alpha_G^*$  is  $\mathbb{G}$ -adapted, Lemma A.1 implies  $\mathbb{E}[z_t(\alpha_G^*)] = \mathbb{E}[\hat{z}_t(\alpha_G^*)|\mathcal{G}_t]$ . The fourth equality follows from the definition of  $z$  and the final equality follows from the definition of  $\alpha_G^*$  as an optimal solution for the left side of (14). If there is no  $\alpha_G^*$  that attains the optimal value (i.e., the supremum is approached but not attained), for any  $\epsilon > 0$  we can find an  $\alpha_G^* \in \mathcal{A}_{\mathbb{G}}$  that is within  $\epsilon$  of the optimal value. The argument above goes through for this  $\alpha_G^*$  except the last line becomes

$$\geq \sup_{\alpha_G \in \mathcal{A}_{\mathbb{G}}} \mathbb{E}[r(\alpha_G) - z(\alpha_G)] - \epsilon.$$

Because  $\epsilon$  can be made arbitrarily small, the desired result still holds.  $\square$

### A.7. Imperfect Information Bounds for the Adaptive Inventory Example

As noted in §3.4, in this case, we would generate demands  $\hat{d}_t^k$  and the corresponding distributions  $\hat{\pi}_t^k$  in the  $k$ th trial of the outer simulation. The inner dual problem is then a stochastic DP that explicitly considers the uncertainty about the ordering costs: the lower-bound cost-to-go function  $\underline{\mathcal{J}}_t^{L,k}(x_t, c_t)$  in the  $k$ th trial can be written recursively as

$$\begin{aligned} \underline{\mathcal{J}}_t^{L,k}(x_t, c_t) &= -c_t x_t + \min_{y_t \geq x_t} \left\{ c_t y_t + \mathbb{E}\left[\underline{\mathcal{J}}_{t+1}^{L,k}(y_t - \hat{d}_t^k, \tilde{c}_{t+1}) - J_{t+1}^{L-1}(y_t - \hat{d}_t^k, \tilde{c}_{t+1}, \hat{\pi}_{t+1}^k)|c_t\right] \right. \\ &\quad \left. + \mathbb{E}\left[f_t(y_t - \tilde{d}_t) + J_{t+1}^{L-1}(y_t - \tilde{d}_t; \tilde{c}_{t+1}, \pi_{t+1}(\hat{\pi}_t^k, \tilde{d}_t))|\hat{\pi}_t^k, c_t\right] \right\} \end{aligned} \quad (32)$$

with the terminal value  $\underline{\mathcal{J}}_T^{L,k}(x_T, c_T) = -c_T x_T$ .

### A.8. Derivation of Equation (22)

Using Ito's lemma and equation (19), we can write the diffusion equation for  $S_\tau = \ln(s_\tau)$  as

$$dS_\tau = (\gamma_\tau - \frac{1}{2}v_\tau)d\tau + \sqrt{v_\tau}dz_\tau^s.$$

Using a discrete-time approximation of this diffusion equation with time steps of length  $\delta$  and taking  $\mathcal{G}_t$  to represent the state of information where all interest rates and volatilities are known,  $S_{t+1}$  is normally distributed with mean and variance:

$$\begin{aligned}\mathbb{E}[S_{t+1}|\mathcal{G}_t] &= S_t + (\gamma_t - \frac{1}{2}v_t)\delta + \rho_{sv}\sqrt{v_t}(v_{t+1} - v_t) \\ \text{Var}[S_{t+1}|\mathcal{G}_t] &= (1 - \rho_{sv}^2)v_t\delta\end{aligned}$$

The stock price  $s_t = \exp(S_t)$  is then log-normally distributed with mean

$$\mathbb{E}[S_{t+1}|\mathcal{G}_t] = \exp\left(\mathbb{E}[S_{t+1}|\mathcal{G}_t] + \frac{1}{2}\text{Var}[S_{t+1}|\mathcal{G}_t]\right)$$

Equation (22) then follows by substituting the expressions above for the mean and variance of  $S_{t+1}$ .

### A.9. Comparison of Option Pricing Bounds with Haugh and Kogan (2004)

As noted in the text, with the perfect information relaxation the “flattened” version of the inner problem for the option pricing example (see equation (24)) is very close to the form of martingale-based duality considered by Haugh and Kogan (2004), Rogers (2002), and Andersen and Broadie (2004), but there is a subtle difference. Taking expectations over equation (24), our upper bound on the value of a call option is given by:

$$\mathbb{E}[\bar{v}_0] = \mathbb{E}\left[\max\left\{\max_{t \in \{0, \dots, T\}} \{\phi_t(s_t - K) - \mu_t\}, -\mu_T\right\}\right] \quad (33)$$

where  $\mu$  is martingale with  $\mu_0 = 0$ . Haugh and Kogan (2004), Rogers (2002), and Andersen and Broadie (2004) consider general martingales (i.e.,  $\mu_0$  need not be 0) and write the dual upper bound on the option as

$$\mathbb{E}\left[\max_{t \in \{0, \dots, T\}} \{\phi_t h_t - \mu_t\}\right] + \mu_0 \quad (34)$$

where  $h_t$  is the option payoff function. Including the  $\mu_0$  term here allows the use of martingales with non-zero initial values, but this is not a substantive difference in formulations: We could always replace  $\mu_t$  with  $\mu_t - \mu_0$  in (34) and have the same bound but with  $\mu_0 = 0$ . Alternatively, we could add  $\mu_0$  to (33).

The subtle difference centers on the definition of the option payoff function  $h_t$  and how non-exercise is handled. Haugh and Kogan and others require  $h_t$  to be a non-negative function that describes the payoffs of the option if exercised in period  $t$ . For a call option, they take  $h_t = \max\{0, (s_t - K)\}$ . There is a small abuse of terminology here: we do not “exercise” an option to get a 0 payoff. We could, however, perhaps throw away or “burn” an option. The possibility of burning an option before expiration doesn't matter in the primal problem, because we would never burn an option before it expires.

However, the possibility of burning an option may matter in the dual. Compare the maximization problems in (33) and (34) in the case of a call option. Problem (34) allows the DM to burn the call option before expiration in period  $t$  ( $t < T$ ) and receive  $-\mu_t$  or exercise the option and receive  $\phi_t(s_t - K) - \mu_t$ . In (33), the DM can receive  $\phi_t(s_t - K) - \mu_t$ , but cannot receive  $-\mu_t$  alone. In other words, for a call option, we have

$$\max_t \{\phi_t h_t - \mu_t\} = \max_t \{\phi_t \max\{0, (s_t - K)\} - \mu_t\} \geq \max\left\{\max_t \{\phi_t(s_t - K) - \mu_t\}, -\mu_T\right\} \quad (35)$$

and it could be the case that the inequality is strict for some  $\mu_t$ . Thus the two formulations (33) and (34) are slightly different, with (33) providing tighter bounds.

## B. Comparison with Stochastic Programming Duality Results

As mentioned in the introduction, the idea of relaxing the nonanticipativity constraints has been exploited in the stochastic programming (SP) literature. In the SP literature, the nonanticipativity constraints discussed in our paper are sometimes called “implementability” constraints. Here we briefly review this SP formulation and compare it to ours. We also briefly compare our formulation to that of Rogers (2007).

### B.1. Stochastic Programming Duality

Our description of the SP approach follows Shapiro and Ruszczyński (2007, pp. 55-75; hereafter SR) and follows them in focusing on perfect information relaxations. The first main assumption is to assume the actions  $a_t$  are scalars or, more generally, vectors in  $\mathbb{R}^n$ . Let  $\alpha_t(\omega)$  denote the action selected in period  $t$  by policy  $\alpha$  with outcome  $\omega$ . The nonanticipativity constraints require the DM to select the same period- $t$  action in all outcomes that are indistinguishable at time  $t$ . We can write these constraints as  $\alpha_t = \mathbb{E}[\alpha_t | \mathcal{F}_t]$ , so the selected action, viewed as a random variable, is equal to its own expected value conditional on the time- $t$  state of information.

SR then place some assumptions on the reward functions and action spaces and use standard Lagrange duality arguments to “dualize” the nonanticipativity constraint. First, assume the reward functions  $(r_0(a, \omega), \dots, r_T(a, \omega))$  depend on the action selected in period  $t$  but are independent of the actions selected in other periods. Second, assume that each  $r_t(a, \omega)$  is polyhedral (piecewise linear and convex) in  $a_t$  for all  $\omega$ . Third, assume the sequences of actions  $a = (a_0, \dots, a_T)$  are drawn from a  $A \subseteq \mathbb{R}^{n(T+1)}$  defined by a set of linear constraints. Now, let  $\lambda_t$  be the Lagrange multipliers associated with the nonanticipativity constraints requiring  $\alpha_t = \mathbb{E}[\alpha_t | \mathcal{F}_t]$ . Since  $\alpha_t$  is a random variable,  $\lambda_t$  is also a random variable (i.e., a function of  $\omega$ ) and has the same dimensionality as  $\alpha_t$ . SR then show that the Lagrangian dual of the stochastic program can be written as

$$\min_{\{\lambda_t: \mathbb{E}[\lambda_t | \mathcal{F}_t] = 0\}} \mathbb{E} \left[ \max_{a \in A} \left\{ \sum_{t=0}^T r_t(a_t) + \lambda_t a_t \right\} \right]. \quad (36)$$

Standard Lagrange duality arguments ensure that weak and strong duality hold in this framework.

In (36), the  $\lambda_t \alpha_t$  term in the objective function can be viewed as analogous to a linear penalty. The constraint  $\mathbb{E}[\lambda_t | \mathcal{F}_t] = 0$  is equivalent to requiring  $\mathbb{E}[\alpha_t \lambda_t | \mathcal{F}_t] = 0$  for all nonanticipative  $\alpha_t$ , which is analogous to our definition of dual feasible penalties (3), albeit with an equality constraint in place of the inequality. The optimization in (36) is no longer constrained by the nonanticipativity constraint and is analogous to the inner problem for the perfect information relaxation in equation (5) above. This allows us to decompose the inner problem into a series of scenario-specific subproblems for a given set of Lagrange multipliers.

Our formulation of the primal DP problem (1) is more general than the SP formulation in that we do not place any restrictions on the action spaces or reward functions and, in the dual, we do not require the penalties to be linear functions of the actions; weak and strong duality hold without these assumptions. As formulated, our examples do not fit within the SP formulation. The option pricing example has a discrete action space. The inventory model has integer constraints on the order quantities and, if even if we ignore these integer constraints, the penalties we considered in the inventory example are not linear functions of the actions and hence are not consistent with the SP formulation.

### B.2. Linear Programming Formulation of DP Duality

As discussed in §2.2, there are strong connections between our results and standard results in linear programming. In fact, if we allow the use of mixed policies, then our formulation of the primal DP can be viewed as a linear programming problem where the decision variables are mixing probabilities on policies; the objective function and constraints are both linear functions of the mixing probabilities. Applying Lagrange duality arguments like those used in the SP framework in our linear programming formulation of the primal delivers results and penalties like ours. However, as shown in §2.2, we can also use simple, direct arguments to establish the duality results. In this subsection, we describe this linear programming formulation and duality argument. For ease of exposition, we will assume that the set of outcomes  $\Omega$  and actions sequences  $A$  are finite sets and, hence, the set of all policies  $\mathcal{A}$  is finite as well, with  $|\mathcal{A}| = |\Omega|^{|A|}$ .

In our “mixed” version of the primal problem (1), rather than selecting a policy  $\alpha$  that selects an action sequence  $a$  in given outcome  $\omega$  (i.e.,  $\alpha : \Omega \rightarrow A$ ), we instead randomly choose a policy  $\alpha \in \mathcal{A}$

using a probability measure  $\mu$ . A mixed policy  $\mu$  is nonanticipative if its mass is concentrated on the set of nonanticipative policies  $\mathcal{A}_{\mathbb{F}}$ . Let  $M$  and  $M_{\mathbb{F}}$  be the set of all mixed policies and all nonanticipative mixed policies, respectively. Clearly  $M$  is a convex set with extreme points corresponding to degenerate distributions that place all of their mass on a single policy  $\alpha$ . Similarly,  $M_{\mathbb{F}}$  is a convex set with extreme points corresponding to degenerate distributions that place all of their mass on a single nonanticipative policy  $\alpha$ .

The mixed version of the original primal (1) can be written as

$$\max_{\mu \in M_{\mathbb{F}}} \mathbb{E}[\mu' \rho], \quad (37)$$

where  $\rho(\omega) = (r(\alpha(\omega), \omega))_{\alpha \in \mathcal{A}}$  is a random vector describing the rewards for each policy  $\alpha$ . (Note  $\rho : \Omega \rightarrow \mathbb{R}^{|\mathcal{A}|}$ .) The inner product  $\mu' \rho$  is a random variable (a function of  $\omega$ ) that represents the expected rewards associated with the mixed policy  $\mu$ , with the expectations taken over the mixture of policies, not the outcomes  $\omega$ . Note the objective function is linear in  $\mu$  and the constraint set  $M_{\mathbb{F}}$  is convex in  $\mu$ ; thus the optimal value will be attained at an extreme point of  $M_{\mathbb{F}}$ . As noted above, the extreme points of  $M_{\mathbb{F}}$  correspond to the degenerate mixed policies that place all of their mass on a single nonanticipative policy  $\alpha$ . Thus the optimal value for the mixed primal (37) will match that of the original primal (1) and each optimal solution for the mixed primal will correspond to a nonanticipative policy that is optimal for the original problem or perhaps a mixture of nonanticipative policies, each of which is optimal for the original problem.

Next we develop a linear equality based representation of the nonanticipativity constraint. First note that we can define nonanticipativity for non-mixed policies using an indicator function  $1_{a_0, \dots, a_t}(a)$  on  $A$  that takes on the value 1 if the first  $t$  elements of  $a$  match  $a_0, \dots, a_t$  and is zero otherwise;  $1_{a_0, \dots, a_t}(\alpha)$  is a random variable (a function of  $\omega$ ) that takes on 1 when  $\alpha$  selects the sequence  $(a_0, \dots, a_t)$  and is zero otherwise. A policy  $\alpha$  is nonanticipative if and only if

$$1_{a_0, \dots, a_t}(\alpha) = \mathbb{E}[1_{a_0, \dots, a_t}(\alpha) | \mathcal{F}_t] \quad \text{for all } t \text{ and } a_0, \dots, a_t. \quad (38)$$

Note that both sides of (38) are random variables and the equality constraints must hold for every  $\omega$ . We now generalize this idea to mixed policies. Let  $\mu_t(a_0, \dots, a_t; \omega)$  denote the probability of choosing the action (sub)sequence  $(a_0, \dots, a_t)$  given outcome  $\omega$  and mixed policy  $\mu$ . This probability  $\mu_t(a_0, \dots, a_t; \omega)$  can be calculated as the inner product  $\mu' \mathbb{1}_{a_0, \dots, a_t}(\omega)$  where  $\mathbb{1}_{a_0, \dots, a_t}(\omega) = (\mathbb{1}_{a_0, \dots, a_t}(\alpha(\omega)))_{\alpha \in \mathcal{A}}$ . Here  $\mathbb{1}_{a_0, \dots, a_t}(\omega)$  is a random vector with entries noting whether policy  $\alpha$  matches the specified action subsequence for the given outcome  $\omega$ . (Note  $\mathbb{1}_{a_0, \dots, a_t} : \Omega \rightarrow \{0, 1\}^{|\mathcal{A}|}$ .) Suppressing the outcome  $\omega$ , we can view  $\mu_t(a_0, \dots, a_t) = \mu' \mathbb{1}_{a_0, \dots, a_t}$  as a random variable. Using this, we can write a linear constraint that requires the probability of choosing an action sequence  $(a_0, \dots, a_t)$  under  $\mu$  to be  $\mathcal{F}_t$ -measurable:

$$\mu_t(a_0, \dots, a_t) = \mathbb{E}[\mu_t(a_0, \dots, a_t) | \mathcal{F}_t] \quad \text{for all } t \text{ and } (a_0, \dots, a_t). \quad (39)$$

If this condition is satisfied, we can build a “decision tree” to describe the expected payoffs of the problem with well-defined probabilities for each decision node that are conditioned on the period- $t$  state of information  $\mathcal{F}_t$  and the prior history of actions. The conditional probabilities for period- $t$  are given by  $\mu_t(a_0, \dots, a_t) / \mu_{t-1}(a_0, \dots, a_{t-1})$  and (39) ensures that these conditional probabilities are measurable with respect to  $\mathcal{F}_t$ , so that they depend only on the outcomes of uncertainties that have already been resolved (e.g., on chance nodes that appear before the decisions in the decision tree).

We can now consider an alternative version of the mixed primal (37) with the linear constraint (39) replacing the nonanticipativity constraint ( $\mu \in M_{\mathbb{F}}$ ):

$$\begin{aligned} \max_{\mu \in M} & \mathbb{E}[\mu' \rho] \\ \text{s.t.} & \mu_t(a_0, \dots, a_t) = \mathbb{E}[\mu_t(a_0, \dots, a_t) | \mathcal{F}_t] \quad \text{for all } t \text{ and } (a_0, \dots, a_t). \end{aligned} \quad (40)$$

It is straightforward to show that any nonanticipative mixture satisfies (39): the nonanticipative mixtures assign positive probability only to policies that satisfy (38) and thus the nonanticipative mixtures must satisfy (39). The constraint (39) may also be satisfied by mixed policies  $\mu$  that are not nonanticipative, so (40) is a relaxation of (37). However, for any mixed policy  $\mu$  satisfying (39), we can construct a nonanticipative mixed policy that is “behaviorally equivalent” to  $\mu$  in that it leads to the same joint probability distribution



on action sequences and outcomes ( $A \times \Omega$ ) and thus leads to the same expected rewards; this follows from a famous result in game theory known as Kuhn's Theorem (see, e.g., Fudenberg and Tirole, 1991).<sup>2</sup> Given this, replacing the nonanticipativity constraint in (37) with the relaxed constraint (39) does not improve the optimal value and we can construct a behaviorally equivalent nonanticipative mixed policy corresponding to any solution to (40).

Having established that the original primal (1) and linear mixed primal (40) have equal optimal values and corresponding solutions, we now proceed to consider the Lagrangian dual of the mixed primal (40) by relaxing the constraints forcing the mixed policies  $\mu$  to satisfy the linear constraints (39). The dual function can be written

$$g(w_t) = \max_{\mu \in M} \mathbb{E} \left[ \rho' \mu - \sum_{t=0}^T w'_t (\mu_t - \mathbb{E} [\mu_t | \mathcal{F}_t]) \right]. \quad (41)$$

The Lagrange multipliers associated with the constraints (39) are given by a stochastic process  $w_t(\omega) \in \mathbb{R}^{|A_t|}$  where  $A_t$  is the set of all possible subsequences of actions  $(a_0, \dots, a_t)$ . (Note that these are the usual Lagrange multipliers divided by  $\mathbb{P}(\{\omega\})$ , so we can bring the Lagrange multipliers inside the expectation in (41).)

Recall that for any  $\mathcal{F}$ -measurable random variable  $X$  and any random variable  $Y$ , we have  $X \mathbb{E}[Y | \mathcal{F}] = \mathbb{E}[XY | \mathcal{F}]$ . Noting this several times and using iterated expectations, we observe:

$$\begin{aligned} \mathbb{E} [w'_t (\mu_t - \mathbb{E} [\mu_t | \mathcal{F}_t])] &= \mathbb{E} [\mathbb{E} [w'_t (\mu_t - \mathbb{E} [\mu_t | \mathcal{F}_t]) | \mathcal{F}_t]] \\ &= \mathbb{E} [\mathbb{E} [w'_t \mu_t | \mathcal{F}_t] - \mathbb{E} [w'_t \mathbb{E} [\mu_t | \mathcal{F}_t] | \mathcal{F}_t]] \\ &= \mathbb{E} [\mathbb{E} [w'_t \mu_t | \mathcal{F}_t] - \mathbb{E} [w_t | \mathcal{F}_t]' \mathbb{E} [\mu_t | \mathcal{F}_t]] \\ &= \mathbb{E} [\mathbb{E} [w'_t \mu_t | \mathcal{F}_t] - \mathbb{E} [\mathbb{E} [w_t | \mathcal{F}_t]' \mu_t | \mathcal{F}_t]] \\ &= \mathbb{E} [\mathbb{E} [\mu'_t (w_t - \mathbb{E} [w_t | \mathcal{F}_t]) | \mathcal{F}_t]] \\ &= \mathbb{E} [\mathbb{E} [\mu'_t z_t | \mathcal{F}_t]] \\ &= \mathbb{E} [\mu'_t z_t], \end{aligned}$$

where  $z_t = w_t - \mathbb{E} [w_t | \mathcal{F}_t]$  and, by construction, we have  $\mathbb{E} [z_t | \mathcal{F}_t] = 0$ . These period- $t$  penalties  $z_t$  are thus constructed like our good penalties as the expectations of the generating functions  $w_t$  that depend on the actions from the first  $t$  periods (and  $\omega$ ). Using this, we can rewrite the dual function (41) as

$$\max_{\mu \in M} \mathbb{E} \left[ \rho' \mu - \sum_{t=0}^T z'_t \mu_t \right]. \quad (42)$$

Any dual feasible  $z_t$ , that is, any  $z_t$  satisfying  $\mathbb{E} [z_t | \mathcal{F}_t] = 0$  (or any set of generating functions  $w_t$  that depend on the first  $t$ -periods actions and  $\omega$ ) will generate an upper bound on the mixed primal problem (40) and hence the original primal (1). Strong Lagrangian duality implies that there exists a  $z_t$  such that the bounds are tight.

Recalling that  $\mu_t = \mu' \mathbf{1}_{a_0, \dots, a_t}$ , we see that this Lagrangian is linear in  $\mu$  and obtains the maximum in (42) at an extreme point of  $M$  which concentrates all of its mass at a single policy  $\alpha$  in  $\mathcal{A}$ . Thus (42) can be written as

$$\max_{\alpha \in \mathcal{A}} \mathbb{E} \left[ r(\alpha) - \sum_{t=0}^T z_t(\alpha) \right]. \quad (43)$$

With no constraints on the policy chosen, we can decompose this into a series of outcome-specific optimization problems where we choose the action  $a$  for each  $\omega$  and rewrite (43) as

$$\mathbb{E} \left[ \max_{a \in A} \left\{ r(a) - \sum_{t=0}^T z_t(a) \right\} \right]. \quad (44)$$

---

<sup>2</sup>Fudenberg, Drew and Tirole, Jean (1991). *Game Theory*. The MIT Press, Cambridge, Massachusetts.

Thus, for any dual feasible penalty  $(z_0, \dots, z_T)$  (or Lagrange multipliers/generating functions  $(w_0, \dots, w_T)$ ), (44) generates an upper bound on the original primal (1) and there exist a penalty/generating function that leads to a tight upper bound. This is exactly the perfect information bound given in equation (5) above.

### B.3. Relationship to Rogers (2007)

Rogers (2007) recently independently proposed an extension of his dual approach to option pricing (Rogers 2002) to Markov decision processes. He considers only the perfect information relaxations and assumes the DP has a Markovian structure with a period- $t$  state variable  $X_t$  that, in our notation, can be viewed as a function of  $\omega$  and the action vector  $a$ . Rogers considers period- $t$  “penalties” of the form  $E[h_{t+1}(X_{t+1})|\mathcal{F}_t] - h_{t+1}(X_{t+1})$ . These penalties are similar to those generated by our Proposition 2.2 except his generating functions  $h_t$  depend on the state  $X_t$  alone whereas our generating functions  $w_t$  can depend on both the outcome  $\omega$  and actions  $a$ .

Rogers shows that weak and strong duality holds with penalties of this form; strong duality is established by taking  $h_t(X_t)$  to be the dynamic programming value function. Rogers provides some ideas and results towards constructing an algorithm for approximately solving a Markov decision problem, but does not consider any specific applications of the approach or numerical examples. Our approach is more general than Rogers in that we do not require the DP to have a Markov structure, we consider imperfect as well as perfect information relaxations, we consider a larger class of penalties, and we present many additional results (e.g., complementary slackness, properties of penalties and relaxations) and some specific examples.

## C. Further Details on the Adaptive Inventory Example

This appendix provides the detailed assumptions for the inventory example that were omitted from the main body of the paper. The seven different state transition matrices are described in Table 5 and the four different prior distributions are shown in Table 6. Figure 2 shows the transition probabilities for the Markov chain for the ordering costs  $c_t$ . Figure 3 shows the lower bounds generated by using the perfect information relaxation with the zero-, one- and two-period look-ahead penalties. The format is the same as the “aquarium plot” of Figure 1. Tables 7 and 8 provide the values and mean standard errors underlying Figure 1 and Figure 3. Table 9 provides the results for the modified myopic policy discussed in §3.6.

Table 5: Transition Probability Matrices.

		Demand distribution ( $\delta_t$ )		
		Low	Medium	High
Stable, Pos. Corr.	Low	0.8	0.1	0.1
	Medium	0.1	0.8	0.1
	High	0.1	0.1	0.8
Stable, Neg. Corr.	Low	0.2	0.4	0.4
	Medium	0.4	0.2	0.4
	High	0.4	0.4	0.2
Stable, Zero Corr.	Low	0.333	0.333	0.333
	Medium	0.333	0.333	0.333
	High	0.333	0.333	0.333
Upward, Slow	Low	0.8	0.1	0.1
	Medium	0.0	0.8	0.2
	High	0.0	0.0	1.0
Upward, Fast	Low	0.2	0.4	0.4
	Medium	0.0	0.2	0.8
	High	0.0	0.0	1.0
Downward, Slow	Low	1.0	0.0	0.0
	Medium	0.2	0.8	0.0
	High	0.1	0.1	0.8
Downward, Fast	Low	1.0	0.0	0.0
	Medium	0.8	0.2	0.0
	High	0.4	0.4	0.2

Table 6: Priors.

	Demand distribution ( $\delta_t$ )		
	Low	Medium	High
H: High	0.10	0.30	0.60
U: Uniform	0.33	0.34	0.33
M: Medium	0.10	0.80	0.10
L: Low	0.60	0.30	0.10

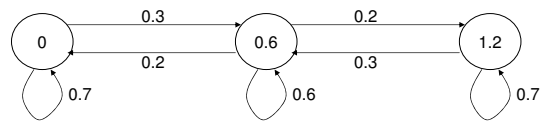


Figure 2: Dynamics of  $c_t$ .

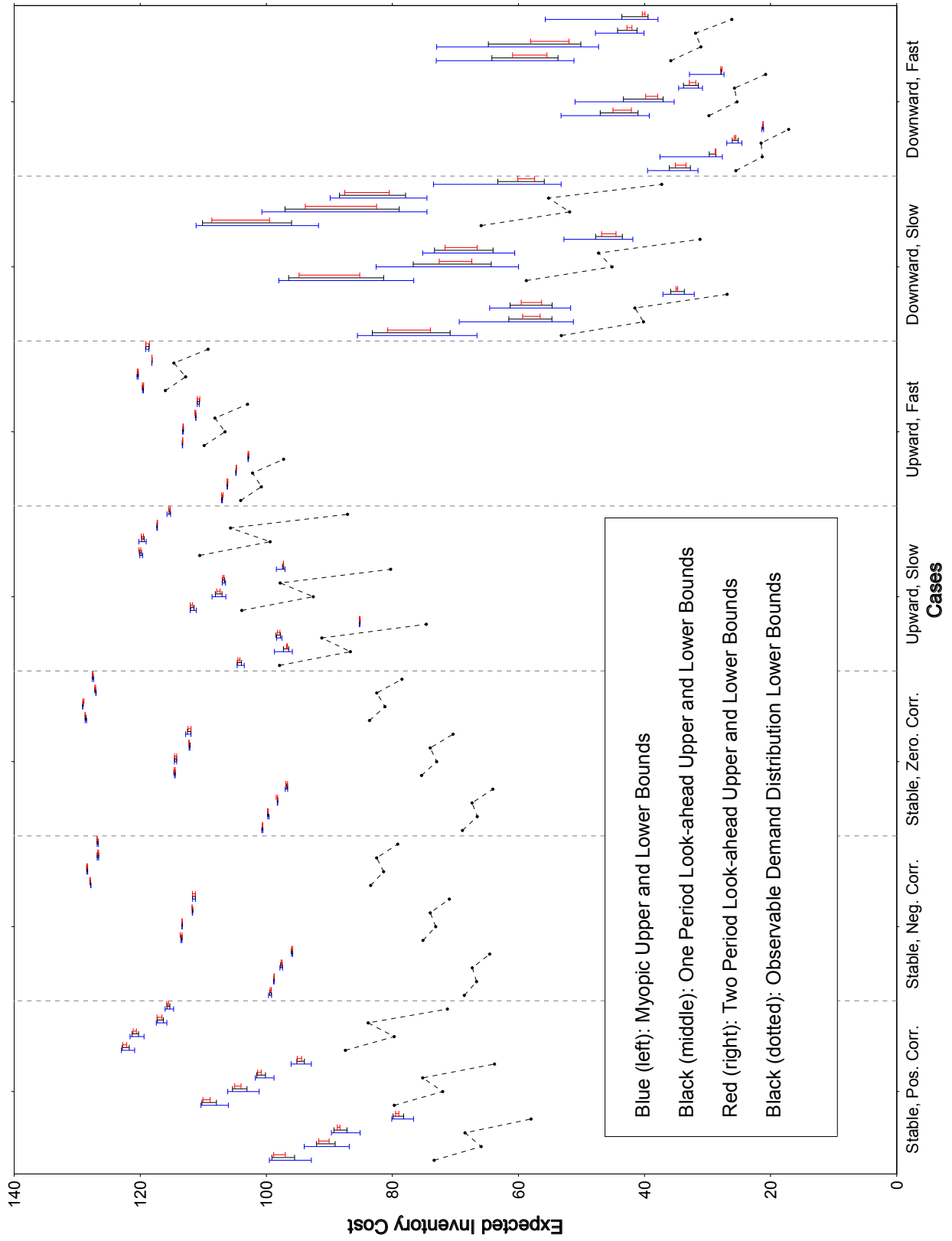


Figure 3: Upper and lower bounds with the imperfect information relaxation.





