

Sketch a theory of nonlinear partial information min-max control

Pierre Bernhard*
INRIA Sophia Antipolis, France

August 1993

Abstract

We sketch a dynamic programming type of theory for both continuous time and discrete time non linear partial information min-max control, using the “cost to come” function as the informational state. We use this theory to derive conditions under which a certainty equivalence principle holds. The condition derived is less powerful than what was known from direct investigation in the continuous time case, but more powerful in the discrete time case.

1 Introduction

We investigate the most classical partial information min max control problem. The set up is as in classical stochastic control, but without probabilistic hypotheses on the disturbances. Instead of minimizing the expected value of a cost criterion, we minimize its supremum with respect to the disturbances.

According to Caratheodory [8], the “cost to come” function has been introduced in the calculus of variations by Hamilton. Carathedory calls it “Hamilton’s principal function”. The conditional cost to come has been used in [3] as an auxiliary tool to solve in a recursive fashion the “auxiliary problem”. It lead to the second (filtering) Riccati equation. Its central role

*The author acknowledges the help of Alain Rapaport in preparing this article

in min-max control has been first stressed, as far as we know, by Didinsky, and used extensively in [9]. John Barras [2] has shown, using tools of large deviation theory in the spirit of [12], that it should be considered as the informational state. The idea of the present derivation is to apply Quadrat's morphism [12],[1] p 442, i.e. to mimic stochastic control, simply replacing mathematical expectations by max operations.

In the continuous time case, problems of differentiability naturally lead to an attempt at using the Fenchel dual of the cost to come, in the spirit of [10], as a possibly more regular informational state, leading to a suboptimal controller if the cost to come is not concave.

2 General framework

2.1 The problem

2.1.1 The dynamics

The general problem considered is as follows. A disturbed control system is given as

$$\dot{x} = f_t(x, u, w), \quad x(0) = x_0, \quad (1)$$

or, in the discrete time case,

$$x_{t+1} = f_t(x, u, w), \quad x(0) = x_0. \quad (2)$$

Here, $t \in [0, T]$ is the time, continuous or discrete, T a fixed horizon. $x \in \mathbb{R}^n$ is the state, $u \in \mathbf{U} \subset \mathbb{R}^m$ is the control. The control function is $u(\cdot) \in \mathbf{U}$. Likewise, $w \in \mathbf{W} \subset \mathbb{R}^l$ is a disturbance, the disturbance function is $w(\cdot) \in \mathbf{W}$.

As the initial state x_0 will also be part of the disturbance, we shall need the notation

$$\omega = (x_0, w(\cdot)) \in \Omega = \mathbb{R}^n \times \mathbf{W}.$$

In the continuous time case, we shall furthermore restrict \mathbf{U} and \mathbf{W} to measurable functions, and the function f will be assumed to be of class C^1 and to satisfy a standard growth hypothesis to insure existence of a unique state trajectory for all $(u(\cdot), \omega) \in \mathbf{U} \times \Omega$.

2.1.2 The information and the strategies

The state x is not directly measured by the controller, but only an output variable $y \in \mathbb{R}^p$:

$$y = h_t(x, w).$$

We let $Y_t = h_t(\mathbb{R}^n, W)$ denote the range of h_t .

As we shall need causal dependencies, we shall use the notation

$$y^t = \{y(\tau) \mid \tau \leq t\},$$

and likewise for all other time functions. Also, we shall use

$$\omega^t = (x_0, w^t) \in \Omega^t$$

In contrast, a subscript will always denote the instantaneous value of a time function, and Ω_t will be a subset of Ω .

In the continuous time case, the control will be allowed to depend in a causal fashion on $y(\cdot)$, and the set of admissible controllers, or strategies, $\mu \in \mathcal{M}$ is that of all causal controllers such that the system obtained by placing

$$u(t) = \mu_t(y) = \mu_t(y^t) \tag{3}$$

in the dynamics generates for all $\omega \in \Omega$ a unique admissible control function $u(\cdot) \in \mathbf{U}$, and thus a unique trajectory.

In the discrete time case, we shall rather let the admissible controllers, again denoted as $\mu \in \mathcal{M}$, be strictly causal :

$$u(t) = u_t = \mu_t(y^{t-1}). \tag{4}$$

(An apparently convenient way to avoid the difference between (3) and (4) would be to put a strict inequality sign in the definition of y^t . However, we believe that this would introduce misleading notations, and we prefer the current convention.)

2.1.3 Criterion and optimality

A criterion or payoff function is given in the classical form, i.e., in the continuous time case as

$$J(u(\cdot), \omega) = M(x(T)) + \int_0^T L_t(x, u, w) dt + N(x_0),$$

and in the discrete time case as

$$J(u(\cdot), \omega) = M(x(T)) + \sum_{t=0}^{T-1} L_t(x, u, w) + N(x_0).$$

As stated above, we are interested in finding a strategy $\mu^* \in \mathcal{M}$ leading to

$$\max_{\omega \in \Omega} J(\mu^*, \omega) = \min_{\mu \in \mathcal{M}} \max_{\omega \in \Omega} J(\mu, \omega). \quad (5)$$

We shall make use of the usual definition of the hamiltonian of the problem:

$$H_t(x, \lambda, u, w) = L_t(x, u, w) + (\lambda, f_t(x, u, w)).$$

In (5) and in the sequel, we write min and max as if they were always reached. We only claim a theory of sufficient conditions. So this is permissible. Specific assumptions could be made to insure the existence of some or all of these extrema.

In the continuous time case the real issue is the regularity of the three value functions U , V , and W below. And we do not tackle this issue.

In the discrete time case things are simpler. One might assume continuity of the functions defining the problem and compactness of the set of initial states $X_0 = \{x \in \mathbb{R}^n \mid N(x) > -\infty\}$. But this would not account for the linear quadratic case. One might use convexity concavity hypotheses, but this would not account for most nonlinear examples. So we chose to stick with this incomplete formulation.

In none of the two cases do we pretend to give a finished theory, as the title of the paper implies.

We shall write

$$\min_u, \quad \max_w, \quad \max_x, \quad \text{for } \min_{u \in \mathcal{U}}, \quad \max_{w \in \mathcal{W}}, \quad \max_{x \in \mathbb{R}^n}$$

respectively.

2.2 Conditional cost to come

As stated above, we shall use the “conditional cost to come” function that we now define.

Let u^τ and y^τ be fixed. Introduce the subset Ω_τ of disturbances that are *compatible* with these data

$$\Omega_\tau[u^\tau, y^\tau] = \{\omega \in \Omega \mid \forall t \leq \tau, h_t(x_t, w_t) = y_t\}. \quad (6)$$

Here, of course, $x(\cdot)$ is the state trajectory generated by $u(\cdot)$ and ω . In this set of *conditional disturbances*, we define the subset $\Omega_\tau(\xi)$ of those that furthermore drive the state to ξ , in the continuous time case at time τ :

$$\Omega_\tau[u^\tau, y^\tau](\xi) = \{\omega \in \Omega_\tau[u^\tau, y^\tau] \mid x(\tau) = \xi\}, \quad (7)$$

and in the discrete time case, at time $\tau + 1$:

$$\Omega_\tau[u^\tau, y^\tau](\xi) = \{\omega \in \Omega_\tau[u^\tau, y^\tau] \mid x(\tau + 1) = \xi\}. \quad (8)$$

We shall oftentimes omit the arguments in square brackets above, as there is no ambiguity in doing so. But it should always be remembered that Ω_t and $\Omega_t(x)$, for instance, depend on u^t and y^t , and thus also the cost to come function defined with them.

The conditional cost to come function $W[u^\tau, y^\tau]$ is now naturally defined.

In the continuous time case :

$$W_\tau(\xi) = \max_{\omega \in \Omega_\tau(\xi)} \left[\int_0^\tau L_t(x, u, w) dt + N(x_0) \right].$$

It should be understood that

$$\forall \xi \mid \Omega_t(\xi) = \emptyset, \quad W_t(\xi) = -\infty.$$

We shall make use of the following hypothesis:

Hypothesis H1 There exists, for all t , u^t , y^t , a well defined cost to go function $W_t(\cdot)$ defined over \mathbb{R}^n into $\mathbb{R} \cup \{-\infty\}$. In the continuous time case, we further assume that it is differentiable with respect to t .

We may have the following further hypothesis:

Hypothesis H1a The function W is of class C^1 .

If W is of class C^1 , it satisfies the forward dynamic programming equation: (remember that u^t and y^t are fixed)

$$\forall x, \quad W_0(x) = N(x), \quad (9)$$

$$\forall(t, x), \quad \frac{\partial W_t(x)}{\partial t} = \max_{w|h_t(x,w)=y_t} H_t(x, -\frac{\partial W_t(x)}{\partial x}, u_t, w) = F_t[W_t, u_t, y_t](x). \quad (10)$$

Let $w = \tilde{\psi}_t(W_t, x, u_t, y_t)$ be a maximizing w above.

In the discrete time case, we have similarly

$$W_\tau(\xi) = \max_{\omega \in \Omega_{\tau-1}(\xi)} \left[\sum_{t=0}^{\tau-1} L_t(x, u, w) + N(x_0) \right],$$

(again, the maximum over an empty set must be taken as $-\infty$.) We use the same hypothesis H1, and the discrete forward dynamic programming equation. We introduce to write it in a readable fashion the following set (again, u^t and y^t are fixed)

$$Z_t(x, u, y) = \{(\xi, v) \in \mathbb{R}^n \times \mathbb{W} \mid f_t(\xi, u, v) = x, \quad h_t(\xi, v) = y\},$$

and it reads

$$W_{t+1}(x) = \max_{(\xi, v) \in Z_t(x, u_t, y_t)} [W_t(\xi) + L_t(\xi, u_t, v)] = F_t[W_t, u_t, y_t] \quad (11)$$

together with (9). Let $(\xi, v) = \zeta_t(W_t, x, u_t, y_t)$ designate a maximizing pair above.

2.3 Cost to go

We shall refer to state feedbacks $u = \phi_t(x)$ and discriminating feedbacks $w = \psi_t(x, u)$. We assume that a consistent framework as been defined (see for instance [4] or [5]) with admissible sets of such feedbacks, within which we have the following property:

Hypothesis H2 The full information game has a unique upper saddle-point (ϕ^*, ψ^*) leading to an upper value denoted $V_t(x)$, of class C^1 in the continuous time case.

Under hypothesis H2, V satisfies Isaacs equation :

$$\forall x, \quad V_T(x) = M(x) \quad (12)$$

and, in the continuous time case,

$$\forall(x, t), \quad \frac{\partial V_t(x)}{\partial t} + \min_{u \in \mathbb{U}} \max_{w \in \mathbb{W}} H_t(x, \frac{\partial V_t(x)}{\partial x}, u, w) = 0 \quad (13)$$

or, in the discrete time case,

$$V_t(x) = \min_{u \in \mathbf{U}} \max_{w \in \mathbf{W}} [V_{t+1}(f_t(x, u, w)) + L_t(x, u, w)] . \quad (14)$$

It should be pointed out that, because we have introduced the term $N(x_0)$ in the payoff, the actual (upper) value of the full information game is

$$\min_{\phi} \max_{\omega} J(\phi, \omega) = \Gamma_0 = \max_x [V_0(x) + N(x)] .$$

3 The continuous time problem

3.1 Dynamic programming

Notice that equations (9) and (10) behave as a filter : (10) can, in principle, be integrated forward in real time from $W_0 = N$, with u_t and y_t as forcing terms in the right hand side. Let \mathcal{W} be the space of all functions from \mathbb{R}^n into \mathbb{R} that can be reached by W_t for all t and all $(u(\cdot), \omega) \in \mathbf{U} \times \Omega$.

Let U_t be a function from \mathcal{W} into \mathbb{R} . Assume it has a derivative with respect to t , denoted $\partial U_t / \partial t$, and a Gateaux derivative with respect to its argument $W \in \mathcal{W}$, denoted DU_t . Actually, it suffices that the directional derivative $DU_t(W).dW$ be defined for all $(W, dW) \in \mathcal{W} \times \mathcal{W}$. The dynamic programming equation is as follows.

$$\forall t, \forall W \in \mathcal{W}, \quad \frac{\partial U_t(W)}{\partial t} + \min_{u \in \mathbf{U}} \max_{y \in \mathbf{Y}_t} DU_t(W).F_t[W, u, y] = 0, \quad (15)$$

where F_t is defined in (10), and the terminal condition

$$\forall W \in \mathcal{W}, \quad U_T(W) = \max_x [M(x) + W(x)]. \quad (16)$$

Notice first that if such a U exists, the minimization over u in (15) defines a control u depending only on t and W . Let $u = \hat{\mu}_t(W)$ be such a minimizing control. Likewise, let $y = \hat{\eta}_t(W, u)$ be a maximizing y . We shall further need the notation $\hat{\psi}_t(W, x, u) = \psi_t(W, x, u, \hat{\eta}(W, u))$.

We can state the following result.

Theorem 3.1 (Main equation) *Assume hypotheses H1 and H1a hold. If there exists a function U from $\mathbb{R} \times \mathcal{W}$ into \mathbb{R} satisfying the main equation*

(15)(16), and corresponding strategies $u_t = \hat{\mu}_t(W_t)$ and $w_t = \hat{\psi}_t(W_t, x, u)$ well defined and admissible when used with (9)(10), then $\hat{\mu}$ is an optimal strategy for the problem at hand, and the optimal cost is $U_0(N)$

Proof As a matter of fact, everything has been cast into the hypotheses, so that the proof is very simple, and similar to that of the classical Isaacs' equation.

Assume first that we play according to $u_t = \hat{\mu}_t(W_t)$. Since, as we have pointed out, W_t only depends on u^t and y^t , this is a causal strategy. Pick a fixed $\bar{\omega} \in \Omega$. Together with $\hat{\mu}$ they generate trajectories $\{x_t\}$, $\{y_t\}$, and $\{W_t\}$ evolving according to (10). By (15), we know that

$$\frac{\partial U_t(W_t)}{\partial t} + DU_t(W_t) \cdot F_t[W_t, \hat{\mu}_t(W_t), y_t] \leq 0. \quad (17)$$

As a consequence, $U_t(W_t)$ will be a decreasing function of t , so that

$$U_T(W_T) \leq U_0(W_0).$$

Let us use (16) and the definition of W_T as cost to come to rewrite the left hand side of the above inequality

$$U_T(W_T) = \max_x \left[M(x) + \max_{\omega \in \Omega_T(x)} \left(\int_0^T L_t dt + N(x_0) \right) \right].$$

Hence

$$U_T(W_T) = \max_{\omega \in \Omega_T} \left[M(x_T) + \left(\int_0^T L_t dt + N(x_0) \right) \right].$$

Now, Ω_T in the above max is that generated by the sequences u^T , y^T generated by $\bar{\omega}$. Therefore, $\bar{\omega} \in \Omega_T$. Thus, using also (9) we get

$$J(\hat{\mu}, \bar{\omega}) \leq \mathbf{U}_T(W_T) \leq U_0(N).$$

Assume now that w plays according to $w_t = \hat{\psi}_t(W_t, x, u)$. On the one hand, because for every (W, x, u, y) , $h_t(x, \tilde{\psi}(W, x, u, y)) = y$, this strategy will generate as output variable $y_t = \hat{\eta}(W_t, u_t)$. Therefore, together with $\hat{\mu}$, we have the equality in (17), and thus $U_T(W_T) = U_0(N)$. On the other hand, w being everywhere maximizing in (10), we have, for such a disturbance,

$J(\hat{\mu}, \omega) = M(x_T) + W_T(x_T)$. Finally, by the hypothesis H1a, W_T is nowhere $-\infty$ over \mathbb{R}^n , i.e., all x 's are reached by maximizing trajectories, so that the right choice of x_0 leads to $J(\hat{\mu}, \omega) = U_T(W_T) = U_0(N)$, which is thus the $\max_{\omega} J(\hat{\mu}, \omega)$.

Assume now that any other admissible strategy μ is chosen by u . Let again w use $w_t = \hat{\psi}_t(W_t, x_t, u_t)$. The same arguments as before show that along the trajectory generated,

$$\frac{\partial U_t(W_t)}{\partial t} + DU_t(W_t) \cdot F_t(W_t, u_t, y_t) \geq 0.$$

Therefore $U_T(W_T) \geq U_0(N)$. And again, there is a choice of x_0 such that for the ω thus generated, $J(\mu, \omega) = U_T(W_T) \geq U_0(N)$. This ends the proof. ■

3.2 An example

The present example is the same as in [7], where a direct derivation of the min-max strategy can be found. By several features it does not fit exactly the previous theory. These features are interesting in that they give hints on how to extend the general theory.

The system is of dimension 2. We call the state variables x and y , and y is also the output variable, i.e. the minimizer's measurement. The dynamics are

$$\dot{x} = -cu + b \cos w \quad (18)$$

$$\dot{y} = -a + b \sin w \quad (19)$$

The constants $a > b > c$ are parameters, the control is $u \in [-1, +1]$ and the disturbance is $w \in [0, 2\pi]$, and x_0 on which we have an a priori information $x_0 \in X_0 = [m_0 - \ell_0, m_0 + \ell_0]$. Since y is observed, so is y_0 , which is always positive.

Final time is $T = \min\{t \mid y(t) = 0\}$, and the payoff is

$$J(u(\cdot), \omega) = |x(T)|.$$

Two features at least do not agree with the above theory.

On the one hand, final time is not given, but is a function of $u(\cdot)$ and ω . But we shall see that knowledge of W allows one to decide whether the

game is terminated or not. Therefore it is trivial to extend the above theory by letting (16) hold for any W such that the game be terminated.

On the other hand, there is a constraint on the allowed initial state x_0 . We may remove it by introducing a term N in the cost, defined by

$$N(x) = \begin{cases} 0 & \text{if } x \in \mathbf{X}_0, \\ -\infty & \text{otherwise.} \end{cases}$$

But now, N is no longer of class C^1 nor will W be. Therefore we cannot use (10), and we shall adapt slightly the theory.

Because of the simple form of the problem (and we shall come back to that point later), we see from its definition that the function W_t will always be of the form

$$W_t(x, y) = \begin{cases} 0 & \text{if } x \in \mathbf{X}_t \text{ and } y = y_t, \\ -\infty & \text{otherwise.} \end{cases}$$

where \mathbf{X}_t is the set of possible x_t 's given the observations up to time t . Given the observation process, this set is always of the form $\mathbf{X}_t = [m_t - \ell_t, m_t + \ell_t]$. And recall that the second component of the state, y_t , is known.

Therefore, the space \mathcal{W} is three dimensional, and $W \in \mathcal{W}$ is entirely characterized by the three variables m , ℓ , and y . It is straightforward to see that the filter is now given by

$$\begin{aligned} \dot{m} &= -cu \\ \dot{\ell} &= \sqrt{b^2 - (\dot{y} + a)^2} \end{aligned}$$

remembering that

$$\dot{y} = b \sin w - a$$

and the main equation now reads

$$U_T(m, \ell, 0) = |m| + \ell$$

and

$$\frac{\partial U_t}{\partial t} + \min_u \max_w [-cu D_m U_t + b(D_\ell U_t |\cos w| + D_y U_t \sin w) - a D_y U_t] = 0.$$

Because of the symmetry with respect to the plane $m = 0$, we may look only at the half space $m \geq 0$. As a matter of fact, we may use a classical

method of characteristics to solve this classical Isaacs equation. We see that the problem being stationary, $\partial U_t / \partial t = 0$, the minimaximizing controls are given by

$$\begin{aligned}\hat{u} &= \text{sign} D_m U \\ \cos \hat{w} &= \frac{D_\ell U}{\sqrt{(D_\ell U)^2 + (D_y U)^2}} \\ \sin \hat{w} &= \frac{D_y U}{\sqrt{(D_\ell U)^2 + (D_y U)^2}}\end{aligned}$$

The adjoint equations are

$$D_m \dot{U} = D_\ell \dot{U} = D_y \dot{U} = 0.$$

In the field of trajectories that reach the terminal plane $y = 0$ at $m > 0$, the transversality conditions give

$$D_m U = 1, \quad D_\ell U = 1,$$

and the main equation gives $D_y U$ as the positive root of the equation

$$-c + b\sqrt{1 + D_y U^2} - aD_y U = 0.$$

This root is thus

$$r = \frac{-ac + \sqrt{a^2 c^2 + (a^2 - b^2)(b^2 - c^2)}}{a^2 - b^2} \quad (20)$$

In that field, one has $\hat{u} = 1$, and \hat{w} given by the above formulas, and this is a certainty equivalent controller with $\hat{x} = m + \ell$.

By placing these controls back into the dynamics in m, ℓ, y , one sees however that the field of trajectories leaves a void between those trajectories that reach $y = 0$ at $m = 0$ (with a negative slope in m) and the symmetry plane $m = 0$. That void may be filled with trajectories reaching the ℓ axis. But there U is not differentiable in m , and the transversality conditions no longer hold. One sees that all trajectories reaching the same point on the ℓ axis will lead to the same U_T , therefore have the same U . And they live in a plane parallel to the m axis, so that in that field, $D_m U = 0$. Thus \hat{u}

is arbitrary as long as it drives the state to the ℓ axis without leaving that field. The adjoint $D_y U$ is now the positive root of

$$b\sqrt{1 + D_y U^2} - aD_y U = 0,$$

let us call it

$$s = \frac{b}{\sqrt{a^2 - b^2}}, \quad (21)$$

leading to a different \hat{w} . It can be seen that the slope of the corresponding trajectories is such that this new field can in fact intersect the first one. One should devise a dividing surface, in fact a plane, which is a dispersal plane in the language of differential games, given by the equality of U in both fields.

The above construction synthesizes a Value function given by

$$U(m, \ell, y) = \max\{|m| + \ell + ry, \ell + sy\}$$

where r and s are the constants given by (20) and (21) respectively.

This analysis gives back the same solution as found in [7], with that difference that there we proposed a particular choice of \hat{u} in the singular field, while we see here that u is arbitrary as long as we do not reach the edge of the field. (The dispersal plane.)

3.3 Certainty equivalence

We can derive from the new theory the following certainty equivalence, which is similar to that of [9], but weaker than that of [7].

Define the *auxiliary problem* as follows. Let again t , u^t , and y^t be fixed, and as usual let Ω_t be $\Omega_t[u^t, y^t]$ as defined in (6). Consider the maximization problem

$$\max_{\omega \in \Omega_t} \left[V_t(x(t)) + \int_0^t L_s(x, u, w) ds + N(x_0) \right].$$

Notice that this problem is equivalent to

$$\max_x [V_t(x) + W_t(x)]. \quad (22)$$

In a sense, the auxiliary problem defines a candidate *worst case* disturbance or, in its second form, worst case current state. We shall call it \hat{x}_t .

We have the following fact.

Theorem 3.2 (Certainty equivalence) *Under the hypotheses of theorem 3.1 and hypothesis H2, if furthermore the auxiliary problem has, for every t and every (u^t, y^t) , a maximum leading to a unique worst case current state \hat{x}_t , the optimum controller of theorem 3.1 is a certainty equivalent strategy $u_t = \phi_t^*(\hat{x}_t)$.*

Proof The proof uses the following fact

Lemma 3.1 *Under the assumptions of the theorem, the function*

$$U_t(W) = \max_x [V_t(x) + W(x)] \quad (23)$$

is a solution of the main equation.

Proof of the lemma Notice first that (23) defines a function that clearly satisfies (16). The hypothesis says that for every $W \in \mathcal{W}$, the max in (23) is reached at a unique \hat{x} . Let us calculate a directional derivative of U_t as given by (23). Let $dW \in \mathcal{W}$, and

$$U_t(W + \theta dW) = \max_x [V_t(x) + W(x) + \theta dW(x)].$$

The maximum in x being reached at a unique point \hat{x}_t , it follows from Danskin's theorem that

$$\frac{\partial}{\partial \theta} U_t(W + \theta dW)|_{\theta=0} = dW(\hat{x}_t).$$

Therefore, U_t has a Gateaux derivative given by

$$DU_t(W).dW = dW(\hat{x}_t).$$

We also use Danskin's theorem to calculate the partial derivative

$$\frac{\partial U_t(W)}{\partial t} = \frac{\partial V_t(\hat{x}_t)}{\partial t}.$$

Placing this back into the main equation (15) with F_t given by (10) yields

$$\frac{\partial V_t(\hat{x}_t)}{\partial t} + \min_{u \in \mathcal{U}} \max_y \left[\max_{w|y} H_t \left(\hat{x}_t, -\frac{\partial W_t}{\partial x}(\hat{x}_t), u, w \right) \right] = 0,$$

i.e.

$$\min_{u \in \mathbf{U}} \max_{w \in \mathbf{W}} H_t \left(\hat{x}_t, -\frac{\partial W_t}{\partial x}(\hat{x}_t), u, w \right) = 0.$$

Notice that since \hat{x}_t maximizes in (23) and both V_t and W_t are assumed C^1 (both are defined over \mathbb{R}^n),

$$\frac{\partial W_t(\hat{x}_t)}{\partial t} = -\frac{\partial V_t(\hat{x}_t)}{\partial t},$$

and thus the previous equation is just Isaacs equation (13). This ends the proof of the lemma.

Let us turn back to the proof of the theorem. The last two equalities show that indeed the minimizing $u_t = \hat{\mu}_t(W_t)$ is just $u_t = \phi^*(\hat{x}_t)$. We see also, using (23) and (9), that in that case, the minimax value of the criterion is

$$U_0(N) = \max_x [V_0(x) + N(x)] = \Gamma_0$$

as in the full information game. ■

Remark If one only assumes uniqueness of the certainty equivalent control $\phi_t^*(\hat{x}_t)$ rather than that of \hat{x}_t itself, the above derivation can easily be extended to show that the main equation is still satisfied with the partial derivative in time replaced by a right derivative, and the Gateaux derivative replaced by a directional derivative. The regularity hypotheses on U then have to be adapted in consequence to let one integrate by parts and reach the same conclusion.

3.4 A dual formulation

In the example above, the payoff is purely terminal, and as a result the cost to go W_t is the characteristic function of the set \mathbf{X}_t of conditionally reachable states defined as

$$W_t(x) = \begin{cases} 0 & \text{if } x \in \mathbf{X}_t, \\ -\infty & \text{otherwise.} \end{cases}$$

If \mathbf{X}_t is convex, this is a concave function. Clearly it is not of class C^1 , but its support function may be, and might be used to characterize W .

More generally, if W is concave u.s.c., it is entirely characterized by its Fenchel dual W^* from \mathbb{R}^n into \mathbb{R} defined as

$$W^*(p) = \min_x [(p, x) - W(x)]. \quad (24)$$

More precisely, the bidual

$$W^{**}(x) = \min_p [(p, x) - W^*(p)]$$

is the smallest concave u.s.c. function larger or equal to W , and is thus W whenever the latter is concave u.s.c.

We shall obtain a meaningful result using W^* even in the case where W is not concave.

We need three more hypotheses :

Hypothesis H3 The minimum over x in (24) is reached at a unique point for every $W \in \mathcal{W}$.

Hypothesis H4 The observation process is with *no reappraisal*, i.e. new information coming in as time goes on cannot invalidate a possible *past* disturbance. Thus, if we let

$$\Omega_t^r = \{(x_0, w^r) \mid (x_0, w(\cdot)) \in \Omega_t\},$$

our new hypothesis states that

$$\forall (u(\cdot), \omega), \quad \Omega_t^t = \Omega_T^t.$$

Hypothesis H5 For all x, w on any maximizing trajectory of the auxiliary problem, the jacobian matrix $\partial h_t / \partial w$ is of full row rank.

The last two hypotheses are satisfied, for instance, if for any (t, x) , the map $w \mapsto h_t(x, w)$ is onto. But this is a more restrictive condition than needed. Observe for instance that it is not satisfied in the example above, while it does satisfy the hypothesis H4, and it shall be a simple matter to formulate it in such a way that H5 be met too.

One should also notice that the fact that the observation process be non anticipative implies that

$$\forall (u(\cdot), \omega), \quad \omega \in \Omega_t \Leftrightarrow \omega^t \in \Omega_t^t.$$

We have the following fact:

Lemma 3.2 *Under hypotheses H1, H3, H4, and H5, W^* satisfies the following dual forward dynamic programming equation:*

$$\frac{\partial W_t^*(p)}{\partial t} + \max_{w|y_t} H\left(\frac{\partial W_t^*(p)}{\partial p}, -p, u_t, w\right) = 0$$

where the notation $w \mid y^t$ should be understood as

$$\{w \mid h_t(\frac{\partial W_t^*(p)}{\partial p}, w) = y_t\}.$$

We shall write the above formula for $\partial W_t^*(p)/\partial t$ as $F_t^*[W_t^*, u_t, y_t]$

Proof Replace W_t by its definition in (24):

$$-W_t^*(p) = \max_x \left[(-p, x) + \max_{\omega \in \Omega_t(x)} \left(\int_0^t L_s ds + N(x_0) \right) \right]$$

Recalling the definition (7) of $\Omega_t(x)$, the above expression yields

$$-W_t^*(p) = \max_{\omega \in \Omega_t} \left[(-p, x_t) + \int_0^t L_s ds + N(x_0) \right].$$

Notice that the expression being maximized in the right hand side only depends on ω^t . According to Hypothesis H4 and to the remark above, we may therefore rewrite the preceding as

$$-W_t^*(p) = \max_{\omega \in \Omega_T} \left[(-p, x_t) + \int_0^t L_s ds + N(x_0) \right].$$

According to Danskin's theorem, we may compute directional derivatives in t and p , that will involve the supremum of some expression over all maximizing ω in the max above. But these expressions depend on ω only through x_t , which is assumed to be the same \hat{x}_t for all, and on w_t for the partial derivative in t . We thus find that W_t^* has a right derivative in t given by

$$-\frac{\partial W_t^*}{\partial t} = \sup_{\hat{w}_t} [(-p, f_t) + L_t](\hat{x}_t, u_t, \hat{w}_t)$$

where the sup is to be taken among those \hat{w}_t that belong to maximizing disturbances. Thanks to hypothesis H5, we may apply standard necessary conditions locally at final time to the constrained maximization problem, and conclude that the optimal w 's maximize the hamiltonian under the constraint. Therefore, the above simply reads

$$-\frac{\partial W_t^*}{\partial t} = \max_{w \mid y_t} [(-p, f_t) + L_t](\hat{x}_t, u_t, w).$$

Similarly, Danskin's theorem yields

$$-\frac{\partial W_t^*}{\partial p} = -\hat{x}_t.$$

The last two equalities yield the result. ■

Therefore, we have a filter type of equation to integrate W^* , the initial condition being of course

$$W_0^*(p) = N^*(p).$$

Finally, let \mathcal{W}^* be the space of all possible functions W_t^* .

Now, the main equation may be replaced by the following one. (We use the same letter U for the value function, which is a slight abuse of notations)

$$\forall t, \forall W^* \in \mathcal{W}^*, \quad \frac{\partial U_t(W^*)}{\partial t} + \min_u \max_y [DU(W^*) \cdot F_t^*[W^*, u, y]] = 0. \quad (25)$$

$$U_T(W^*) = \max_x \min_p [(p, x) + M(x) - W^*(p)] \quad (26)$$

Again, call $\hat{\mu}_t(W^*)$ a minimizing control in (25). We have the following result ;

Theorem 3.3 *Under hypotheses H1, H3, H4, and H5, if there exists a function U satisfying (25),(26), and a corresponding control $\hat{\mu}_t(W_t^*)$ generating an admissible strategy, this control guarantees a performance no worse than $U_0(N^*)$, and if W_T is concave, this controller is optimal.*

Proof The proof is completely similar to that of the main theorem, up to the fact that, by playing $u_t = \hat{\mu}_t(W_t^*)$, we get

$$U_T(W_T^*) \leq U_0(W_0^*)$$

with the equality possible for some ω , and the reverse inequality for some ω for any other admissible controller. We end up by noticing that for any control,

$$\begin{aligned} \max_{\omega} J(u(\cdot), \omega) &= \max_x [M(x) + W_T(x)] \leq \max_x [M(x) + W_T^{**}(x)] \\ &= \max_x [M(x) + \min_p [(p, x) - W_T^*(p)]] = U_T(W_T^*) \end{aligned}$$

And if W_T is concave, the only inequality above is an equality.

3.5 The example : dual formulation

To apply the above theory, we convert the measurement process into one that satisfies hypothesis H5 by pretending that our measurement is \dot{y} . The constraint now reads

$$-a + b \sin w = \dot{y}_t .$$

Let us apply the above theory, with (p, q) as the dual variable. (p in the general theory.) It gives, first for W^* :

$$W_0^*(p, q) = pm_0 - |p|l_0 + qy_0 ,$$

and

$$-\frac{\partial W_t^*}{\partial t} = \max_{w|b \sin w = \dot{y}_t} (pcu - pb \cos w + qa - qb \sin w) ,$$

i.e.

$$\frac{\partial W_t^*}{\partial t} = -pcu - |p|\sqrt{b^2 - (\dot{y}_t + a)^2} + q\dot{y}_t .$$

This is a particularly simple P.D.E. which yields

$$W_t^*(p, q) = \alpha p + yq ,$$

with

$$\alpha_0 = m_0 - \text{sign}(p)l_0$$

and

$$\dot{\alpha} = -cu - \text{sign}(p)\sqrt{b^2 - (\dot{y} + a)^2} .$$

This last equation leads to

$$\alpha_t = \mu_t - \text{sign}(p)\lambda_t$$

with

$$\begin{aligned} \lambda_t &= l_0 + \int_0^t \sqrt{b^2 - (\dot{y}_s + a)^2} ds , \\ \mu_t &= m_0 - c \int_0^t u_s ds . \end{aligned}$$

The main equation (25) is now converted into

$$\frac{\partial U_t(W^*)}{\partial t} + \min_u \max_w DU_t(W^*) \cdot [-pcu + |p \cos w| - qa + qb \sin w] = 0 .$$

and if the space \mathcal{W}^* is made up of functions of the form above, time T is defined by $\partial W^*/\partial q = 0$, and the final condition on U is

$$U_T(W^*) = \max_x \min_p [px - \alpha_T p + |x|].$$

It is a simple matter to check that this gives

$$U_T(W^*) = |\mu_T| + \lambda_T.$$

We may now check that the function

$$U_t(W^*) = \max \left\{ \begin{aligned} & \frac{\partial W^*}{\partial p}(-1, 0) + r \frac{\partial W^*}{\partial q}(-1, 0), \\ & \frac{1}{2} \left(\frac{\partial W^*}{\partial p}(-1, 0) - \frac{\partial W^*}{\partial p}(1, 0) \right) + s \frac{\partial W^*}{\partial q}(1, 0), \\ & -\frac{\partial W^*}{\partial p}(1, 0) + r \frac{\partial W^*}{\partial q}(1, 0) \end{aligned} \right\},$$

or, in a simpler form,

$$U_t(W^*) = \max \{ |\mu| + \lambda + ry, \lambda + sy \}$$

is indeed a solution of the main equation. One has to identify the Gateaux derivative of that function and apply it to $F_t^*[W^*, u, y](\cdot)$ as given previously, a simple matter again in this case, as one can see that each of the three terms in the big max above is a linear operator on W^* , and should therefore be applied unchanged on $F_t^*[W^*, u, y](\cdot)$ in the region where it holds.

4 The discrete time problem

4.1 Dynamic programming

The “filter” F_t is now defined by (11), with the same initial condition (9) as in the continuous time case.

The dynamic programming equation now reads :

$$U_t(W) = \min_u \max_y U_{t+1}(F_t[W, u, y]) \tag{27}$$

with the same end condition (16) as in the continuous time case:

$$U_T(W) = \max_x [M(x) + W(x)].$$

If such a function U exists, the above minimization leads to a control $u_t = \hat{\mu}_t(W)$ depending only on W . Notice also that if W_t is computed according to (9)(11), it only depends on past u 's and past y 's, up to time $t-1$. It is therefore possible to choose the minimizing control as $u_t = \hat{\mu}_t(W_t)$. Let us call that strategy $u = \mu^*(y)$. Let again $y = \hat{\eta}_t(W, u)$ be a maximizing y in (27).

One should notice that equations (16) and (27) imply that if $[-\infty]$ is the constant function equal to $-\infty$, then $U_t([-\infty]) = -\infty$. As a consequence, $\hat{\eta}_t(W, u)$ always has an inverse image by h_t , with an x component such that $W_t(x) > -\infty$.

Theorem 4.1 *Under the hypotheses H1 and H2, the above strategy is min-max among the strictly causal strategies. It leads to a performance*

$$\max_{\omega} J(\mu^*, \omega) = U_0(N).$$

Proof Assume the strategy μ^* is used. Pick a fixed $\bar{\omega}$. Together, they generate sequences $\{y_t\}$ and $\{W_t\}$. Along a trajectory, the main equation (27) gives,

$$U_t(W_t) \geq U_{t+1}(W_{t+1}). \quad (28)$$

Therefore, $U_T(W_T) \leq U_0(N)$. As in the continuous time case, this results in

$$J(\mu^*, \bar{\omega}) \leq U_0(N).$$

Consider the sequence $\{\hat{W}_t\}$ generated by (9) and (11) in which we place $u_t = \hat{\mu}_t(\hat{W}_t)$ and $y_t = \hat{\eta}_t(\hat{W}_t, \hat{\mu}_t(\hat{W}_t))$ (and not $y_t = h_t(x_t, w_t)$: at this stage, we have not specified a w_t .) Let \hat{x}_T be such that

$$U_T(\hat{W}_T) = \max_x [M(x) + \hat{W}_T(x)] = M(\hat{x}_T) + \hat{W}_T(\hat{x}_T).$$

From this, construct the backward sequence

$$(\hat{x}_t, \hat{w}_t) = \zeta(\hat{W}_t, \hat{x}_{t+1}, \hat{\mu}_t(\hat{W}_t), \hat{\eta}_t(\hat{W}_t, \hat{\mu}_t(\hat{W}_t)))$$

By construction, this generates a feasible trajectory, generated by $\hat{\mu}$ and a $\hat{\omega} = (\hat{x}_0, \{\hat{w}_t\})$. On the one hand, the output generated by this trajectory is precisely $y_t = \hat{\eta}_t(\hat{W}_t, \hat{\mu}_t(\hat{W}_t))$, which was used to construct the filter. As a consequence, along that trajectory, we shall have an equality in (28), and thus $U_T(\hat{W}_T) = U_0(N)$. On the other hand, because (\hat{x}_t, \hat{w}_t) maximize in (11), we get

$$J(\hat{\mu}, \hat{\omega}) = M(\hat{x}_T) + \hat{W}_T(\hat{x}_T) = U_T(\hat{W}_T) = U_0(N).$$

Thus this is the max in ω of $J(\hat{\mu}, \omega)$.

For any other admissible strategy μ , consider likewise the sequence $\{W_t\}$ generated by $u_t = \mu_t$, $y_t = \hat{\eta}(W_t, \mu_t)$. It will cause $U_{t+1}(W_{t+1}) \geq U_t(W_t)$, hence $U_T(W_T) \geq U_0(N)$, and as above, we may exhibit an ω which yields precisely that sequence $y_t = \hat{\eta}_t(W_t, u_t)$ and $J(\mu, \omega) = U_T(W_T)$. This ends the proof. ■

4.2 Certainty equivalence principle

We define as previously the auxiliary problem in either its extensive form

$$\max_{\omega \in \Omega_t} \left[V_t(x(t)) + \sum_{s=0}^{t-1} L_s(x, u, w) + N(x_0) \right],$$

or in its equivalent form (22). We also let

$$S_t(x, u) = \max_{w \in W} [V_{t+1}(f_t(x, u, w)) + L_t(x, u, w) + W_t(x)].$$

We have the theorem:

Theorem 4.2 *Under hypotheses H1 and H2, if furthermore $\forall (u(\cdot), \omega) \in \mathbf{U} \times \Omega$ and $\forall t$ the function S_t has a saddle-point*

$$\max_x \min_{u \in \mathbf{U}} S_t(x, u) = \min_{u \in \mathbf{U}} \max_x S_t(x, u),$$

then the certainty equivalent controller $u_t = \phi^(\hat{x}_t)$, where \hat{x}_t is a current worst case state in the auxiliary problem, is unique and optimal.*

Proof As in the continuous time case, we show that the function (23) is a solution of the main dynamic programming equation.

Lemma 4.1 *Under the assumptions of the theorem, the function*

$$U_t(W) = \max_x [V_t(x) + W(x)]$$

is a solution of the main equation.

Proof of the lemma Assume the above form for U_t . The left hand side of the main equation (27) can be rewritten using Isaacs' equation (14) to substitute for V_t :

$$U_t(W_t) = \max_x \min_u \max_w [V_{t+1}(f_t(x, u, w)) + L_t(x, u, w) + W_t(x)]. \quad (29)$$

The right hand side requires some more work. We have, assuming the special form for U_{t+1} ,

$$\begin{aligned} & \max_y U_{t+1}(F_t[W_t, u, y]) \\ &= \max_y \max_x \left\{ V_{t+1}(x) + \max_{(\xi, v) \in Z_t(x, u, y)} [W_t(\xi) + L_t(\xi, u, v)] \right\} \\ &= \max_{\xi, v} \{ V_{t+1}(f_t(\xi, u, v)) + W_t(\xi) + L_t(\xi, u, v) \}. \end{aligned}$$

Equivalently, changing the name of the mute variables from (ξ, v) to (x, w) , we obtain

$$\begin{aligned} & \min_u \max_y U_{t+1}(F_t[W_t, u, y]) = \\ & \min_u \max_x \max_w [V_{t+1}(f_t(x, u, w)) + W_t(x) + L_t(x, u, w)] \quad (30) \end{aligned}$$

The comparison of (29) and of (30) yields the result of the lemma.

Unicity of the certainty equivalent control follows from the following :

Lemma 4.2 *Under hypothesis H1, existence of a saddle-point to S implies that the certainty equivalent control $\phi_t^*(\hat{x}_t)$ is unique.*

Proof of the lemma If the current worst case state \hat{x}_t is unique, then since we have assumed uniqueness of the optimal feedback strategy ϕ^* , the result follows. Assume that at some time instant, \hat{x}_t is not unique. Let therefore

\hat{x}^1 and \hat{x}^2 be two states that maximize in (22), and thus also in (29) above. Let $\phi^*(\hat{x}^1) = u^1$ and $\phi^*(\hat{x}^2) = u^2$, and

$$\max_x \min_u S_t(x, u) = S_t(\hat{x}^1, u^1) = S_t(\hat{x}^2, u^2) = \Gamma_t.$$

Since the optimal feedback strategy ϕ^* is unique, it follows that

$$\forall u \neq u^1, \quad S_t(\hat{x}^1, u) \geq \Gamma_t.$$

Thus if $\min_u \max_x S_t = \Gamma_t$, it can only be reached at $u = u^1$. But if $u^1 \neq u^2$, the same uniqueness hypothesis yields

$$S_t(\hat{x}^2, u^1) \geq \Gamma_t.$$

Thus necessarily $u^2 = u^1$. This ends the proof of the lemma.

The proof of the certainty equivalence theorem follows. As a matter of fact, the optimal u is that of (30), but it then coincides with that of (29), which is just $\phi_t^*(\hat{x}_t)$. ■

Remark As far as we know, this is the most general statement to date of a condition under which a discrete-time certainty equivalence principle holds. The first such principle ever seems to have been in [13], but was restricted to the linear quadratic case in its approach. In [6], we gave a proof with convexity concavity hypotheses that hardly hold in any case without at least linear dynamics. In [3], we hinted at the possibility to extend the linear theory of the book to a nonlinear setup. But without a detailed statement. An interesting feature of the condition given here is that it implies uniqueness of the certainty equivalent control rather than of the worst case state.

5 Conclusion

It is not clear that the above theory is of much practical use as such. One of its merits is to give a more general framework within which to understand the certainty equivalence principle, and also to clarify the said principle in the discrete time case. Other particular cases might be interesting to investigate.

We want to make a remark concerning the value functions U_t : they are always max-plus linear (m.p.l.) in the sense of [1]. That is, given two functions W^1 and W^2 from \mathbb{R}^n into \mathbb{R} and given a real number a considered

as the constant function from \mathbb{R}^n into \mathbb{R} equal to a for all x , we have, for all t ,

$$\begin{aligned} U_t(\max\{W^1, W^2\}) &= \max\{U_t(W^1), U_t(W^2)\}, \\ U_t(a + W) &= a + U_t(W). \end{aligned}$$

One should be careful, however, that the \max in the argument of the first equation above is to be taken in the sense of the partial ordering of real functions :

$$W^1 \leq W^2 \Leftrightarrow \forall x \in \mathbb{R}^N, \quad W^1(x) \leq W^2(x),$$

and may not exist. This is *not* to be confused with the function

$$\max\{W^1, W^2\} : x \mapsto \max\{W^1(x), W^2(x)\}.$$

Let us first prove the claim:

Proposition 5.1 *In the continuous time case, under the hypotheses of the main theorem, the function U_t is m.p.l. for all t . In the discrete time case, a function satisfying the main equation (27), (16) is m.p.l.*

Proof We consider the two cases in the reverse order.

In the discrete time case, notice first that the formula (16) for U_T is trivially m.p.l. Now, take the recurrence formula (27) and assume that U_{t+1} is m.p.l. Replacing W by a larger (in the partial ordering of real functions) function \bar{W} in F_t as given by (11) clearly gives, for any (u, y) a larger function $F_t[\bar{W}, u, y]$, and therefore $U_t(\bar{W}) \geq U_t(W)$. Similarly, adding a constant to W in F_t adds that constant to $F_t[W, u, y]$, and by the recurrence hypothesis also to $U_t(W)$.

In the continuous time case, it would be nice to derive the m.p.l. property directly from the functional differential equation (15), as the terminal condition is the same as in the discrete time case, and thus m.p.l. also. This seems less easy to do though. However, under the hypotheses of the main theorem (3.1), the function U_t has an interpretation as

$$U_t(W) = \min_{\mu \in \mathcal{M}(x_t, w(\cdot))} \max \left[M(x_T) + \int_t^T L_s(x_s, u_s, w_s) ds + W(x_t) \right]$$

under the constraint (1). This is again clearly m.p.l., as replacing W by a larger \bar{W} again clearly increases $U_t(W)$, and adding a constant to W adds the same constant to $U_t(W)$. ■

This shows that there is still some unexploited structure in that theory, and thus a hope to exploit it further.

References

- [1] F. Baccelli, G. Cohen, J-P. Quadrat, G-J. Olsder: *Synchronization and Linearity*, Wiley, Chichester, U.K., 1992.
- [2] J. Baras: “Risk Sensitive Control and Dynamic Games for Partially Observed Discrete Time Nonlinear Systems”, workshop on Robust Controller Designs and Differential Games, UCSB, Santa Barbara, California, 1993.
- [3] T. Başar and P. Bernhard: *H_∞ Optimal Control and Related Minimax Design Problems*, Birkhauser, Boston, Mas. 1991.
- [4] P. Bernhard: “Differential Games, closed loop”, *Encyclopaedia of Systems and Control*, Pergamon, London, U.K., 1988.
- [5] P. Bernhard: “Lecture notes on the Isaacs-Breakwell theory”, Summer school on Differential Games, Cagliari, Italy, 1992. To appear.
- [6] P. Bernhard: “A min-max certainty equivalence principle and its applications to continuous time, sampled data, and discrete time H_∞ optimal control”, INRIA Research Report 1347, 1990.
- [7] P. Bernhard and A. Rapaport: “Min-max certainty equivalence principle and differential games”, revised version of a paper presented at the Workshop on Robust Controller Design and Differential Games, UCSB, Santa Barbara, California, 1993. Submitted for publication.
- [8] C. Caratheodory: *Calculus of Variations and Partial Differential Equations of the First Order*, Holden Day, 1967. (Original in German, Teubner, Berlin, 1935)
- [9] G. Didinsky, T. Başar, and P. Bernhard: “Structural Properties of Minimax Policies for a Class of Differential Games Arising in Nonlinear H^∞ -Control”, *Systems and Control Letters*, pp ,1993.
- [10] N.N. Krassovski et A.I. Subbotin: *Jeux Différentiels*, MIR, Moscow, 1977.

- [11] Max Plus: “A linear system theory for systems subject to synchronization and saturation constraints”, *Proceedings of the first European Control Conference*, pp 1022–1033, Grenoble, France, 1991. Hermès, Paris.
- [12] J-P. Quadrat: “Théorèmes asymptotiques en programmation dynamique”, *Compte Rendus de l’Académie des Sciences*, 311: pp 745–748, 1990.
- [13] P. Whittle, “Risk sensitive Linear/Quadratic/Gaussian Control”, *Advances in Applied Probabilities* **13**, pp 764–777, 1981.