

ON BANG - BANG CONTROL POLICIES

Roberto GONZALEZ (*) and Edmundo ROFMAN (**)
Instituto de Matemática "Beppo Levi"
Universidad Nacional de Rosario
ARGENTINA

(Work included in the cooperation program with I.R.I.A. - Rocquencourt - FRANCE)

ABSTRACT: In this paper it is proposed a method for the determination of the optimal distribution of N switching points for a bang-bang control applied to a differential system.

After pointing the necessary conditions to be verified by such switching points it is showed the existence of an optimal policy for a fixed number N of them.

Once characterized these points thru the application of the Pontryagin principle the problem, considered till now in the space of step functions, is put into the L_1 space, in order to show the existence of a minimizing succession of the amplified problem and analyzed its correspondence to an optimal policy.

After reducing the problem into one of optimization on a convex K of \mathbb{R}^n there are added considerations which let us, with the proposed method, obtain the optimal also with a number of switching points n less than the predetermined N .

Now it is proved that the function to optimize is of C^2 class in K and the applied methods are these of the projected gradient and the conjugated gradient conveniently penalized.

Finally, the obtained algorithms are applied in one example: the shut down policy of a nuclear reactor where the optimum is obtained with a finite number of switchings; this number remains constant although increasing values of N are proposed.

§1. STATEMENT OF THE PROBLEM, NECESSARY CONDITIONS OF OPTIMALITY, EXISTENCE OF MINIMUM IN THE CASE OF FIXED NUMBER OF SWITCHING POINTS.

Given the dynamical system governed by the differential equation:

$$(1) \quad \dot{x} = F(t)x + G(t)u \quad x \in \mathbb{R}^n, \quad u \in \mathbb{R}^1$$

with initial condition $x(0) = x_0$

and the cost functional

$$(2) \quad J(u(\cdot)) = \int_0^T l(x(s), u(s), s) ds + g(x(T))$$

we try to find the control function $u(\cdot)$ that minimizes J .

(*) Researcher of the "Consejo de Investigaciones de la Universidad Nacional de Rosario" for the project: "Optimization and Control, Theory and applications".

(**) Director of the above referred project.

The control $u(\cdot)$ belongs to the family \mathcal{U}_{ad} that satisfy the following restrictions.

- a) $u(t) = v_1$ or $u(t) = v_2 \quad \forall t \in [0, T]$
- b) $u(\cdot)$ is a step-function with n switchings
- c) $u(0) = v_1$.

We denote with $\theta_1, \theta_2, \dots, \theta_n$ the switching points and this set with the vector $\theta = (\theta_1, \dots, \theta_n)'$

(3) θ satisfies the restrictions: $0 < \theta_1 < \theta_2 < \dots < \theta_n < T$.

Then, if we fix θ , we know the value of $u(t) \quad \forall t \in [0, T]$ and we can think of $J(u(\cdot))$ as a function $J(\theta)$ of $\theta \in \mathring{\Omega}$, where $\mathring{\Omega} = \{\theta \in \mathbb{R}^n / 0 < \theta_1 < \dots < \theta_n < T\}$. If the minimum of the problem exists, the necessary conditions will be (because θ belongs to an open set):

$$\frac{\partial J}{\partial \theta_1} = 0 \quad ; \quad \dots \quad ; \quad \frac{\partial J}{\partial \theta_n} = 0$$

We can modify the restrictions (3) in the following form:

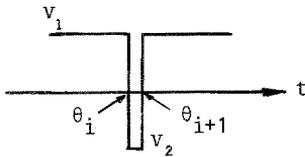
(3') $0 \leq \theta_1 \leq \theta_2 \leq \dots \leq \theta_n \leq T$

and analyze the meaning of a point in the boundary of $\mathring{\Omega}$.

- a) $\theta_1 = 0$ means that the first step has the value $u = v_2$.

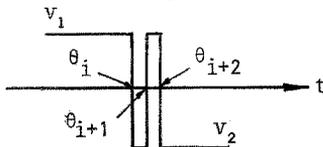
there are two simultaneous commutations that could be eliminated and it remains a new control function with $n-2$ commutations.

- b) $\theta_i = \theta_{i+1}$



there are three simultaneous switchings, we could eliminate two and obtain a new control with $n-2$ commutations.

- c) $\theta_i = \theta_{i+1} = \theta_{i+2}$



- d) $\theta_n = T$ a commutation at the end that could be eliminated and the new control has $n-1$ switchings.

With this meaning we can define the function $J(\theta)$ in the set $\Omega = \{\theta \in \mathbb{R}^n / 0 \leq \theta_1 \leq \theta_2 \leq \dots \leq \theta_n \leq T\}$. In this compact set, under suitable conditions on l, g, F, G , (it will be enough the continuity of l and g , and that F, G be integrable), is $J(\theta)$ continuous; then there is an optimal control in Ω that provides the minimum value of J and has $n' \leq n$ switching points.

§2. THE RELATION BETWEEN THE NECESSARY CONDITION $\frac{\partial J}{\partial \theta_i} = 0$ AND PONTYAGIN'S MAXIMUM PRINCIPLE.

We shall see in the following number that $\frac{\partial J}{\partial \theta_i}$ has the form:

$$\frac{\partial J}{\partial \theta_i} = l(x(\theta_i), u(\theta_i^-), \theta_i) - l(x(\theta_i), u(\theta_i^+), \theta_i) + p(\theta_i)G(\theta_i) [u(\theta_i^+) - u(\theta_i^-)]$$

where $p(t)$ satisfies:

$$\begin{cases} \frac{dp}{dt} = -p F(t) + \frac{\partial l}{\partial x}(x(t), u(t), t) \\ p(T) = -\frac{\partial g}{\partial x} \Big|_{x(T)} \end{cases}$$

$$\text{and } u(\theta_i^-) = \lim_{\substack{\epsilon \rightarrow 0 \\ \epsilon > 0}} u(\theta_i - \epsilon) \quad ; \quad u(\theta_i^+) = \lim_{\substack{\epsilon \rightarrow 0 \\ \epsilon > 0}} u(\theta_i + \epsilon) \quad .$$

If we define:

$$H(x, p, u, t) = p(t)[F(t)x(t) + G(t)u(t)] - l(x(t), u(t), t)$$

we can write:

$$\frac{\partial J}{\partial \theta_i} = H(x(\theta_i), p(\theta_i), u(\theta_i^+), \theta_i) - H(x(\theta_i), p(\theta_i), u(\theta_i^-), \theta_i) \quad .$$

We also know that in the problem $\min_{u \in \mathcal{U}_{ad}} J(u(\cdot))$, where \mathcal{U}_{ad} is the set of step

functions with values v_1, v_2 , if $\bar{u}(\cdot)$ is the optimal control in this set, it must satisfy the maximum principle:

$$H(x(t), p(t), \bar{u}(t), t) = M(x(t), p(t), t) \quad \text{a.e.}$$

where, by definition

$$M(x, p, t) = \max_{u = \begin{cases} u_1 \\ u_2 \end{cases}} H(x, p, u, t)$$

$M(x(t), p(t), t)$ is continuous in t , then:

$$\lim_{\epsilon \rightarrow 0^+} M(x(t+\epsilon), p(t+\epsilon), t) = \lim_{\epsilon \rightarrow 0^-} M(x(t+\epsilon), p(t+\epsilon), t+\epsilon)$$

but

$$\lim_{\epsilon \rightarrow 0^+} M(x(t+\epsilon), p(t+\epsilon), t+\epsilon) = \lim_{\epsilon \rightarrow 0^+} H(x(t+\epsilon), p(t+\epsilon), u(t+\epsilon), t+\epsilon) = H(x(t), p(t), u(t^+), t)$$

and also

$$\lim_{\epsilon \rightarrow 0^-} M(x(t+\epsilon), p(t+\epsilon), t+\epsilon) = \lim_{\epsilon \rightarrow 0^-} H(x(t+\epsilon), p(t+\epsilon), u(t+\epsilon), t+\epsilon) = H(x(t), p(t), u(t^-), t)$$

from where it follows:

$$\frac{\partial J}{\partial \theta_i} = H(x(\theta_i), p(\theta_i), u(\theta_i^+), \theta_i) - H(x(\theta_i), p(\theta_i), u(\theta_i^-), \theta_i) = 0$$

then, the maximum principle implies the necessary conditions of optimality: $\frac{\partial J}{\partial \theta_i} = 0$.

§3. COMPUTATION $\frac{\partial J}{\partial \theta_i}$

The equation of the system's evolution is

$$x(t) = \Phi(t, 0)x_0 + \sum_{k=1}^{j-1} \int_{\theta_k}^{\theta_{k+1}} \Phi(t, s)G(s)u(s)ds + \int_{\theta_j}^t \Phi(t, s)G(s)u(s)ds$$

where $j / \theta_j \leq t \leq \theta_{j+1}$

and we can calculate $\frac{\partial}{\partial \theta_i} x(t)$.

We suppose $G(t)$ is continuous in the interval $[0, T]$

a) $t > \theta_i$

$$\frac{\partial}{\partial \theta_1} x(t) = \Phi(t, \theta_1) G(\theta_1) (v_2 - v_1) (-1)^i .$$

It must be remembered that $\Phi(t, s)$ is the solution of the matrix differential equation

$$\frac{d}{dt} \Phi(t, s) = F(t) \Phi(t, s) \quad \Phi \text{ } v \times v \text{ matrix}$$

with initial condition: $\Phi(s, s) = I$.

b) $t < \theta_1$

$x(t)$ does not depend on θ_1 , then $\frac{\partial}{\partial \theta_1} x(t) = 0$.

In this form, we can say that

$$\begin{cases} \frac{d}{dt} \left(\frac{\partial}{\partial \theta_1} x(t) \right) = F(t) \left(\frac{\partial}{\partial \theta_1} x(t) \right) & t > \theta_1 \\ \text{with initial conditions:} \\ \left. \frac{\partial}{\partial \theta_1} x(t) \right|_{t=\theta_1} = G(\theta_1) (v_2 - v_1) (-1)^i \end{cases}$$

$$J(\theta) = g(x(T)) + \int_0^T l(x(s), u(s), s) dt$$

if l and g are continuously differentiable, we have:

$$\begin{aligned} \frac{\partial J}{\partial \theta_1} &= l(x(\theta_1), u(\theta_1^-), \theta_1) - l(x(\theta_1), u(\theta_1^+), \theta_1) + \\ &+ g'(x(T)) \cdot \frac{\partial}{\partial \theta_1} x(T) + \int_{\theta_1}^T \frac{\partial l}{\partial x}(x(s), u(s), s) \cdot \frac{\partial}{\partial \theta_1} x(s) ds \end{aligned}$$

if we introduce the adjoint vector $p(t)$ that satisfies

$$\begin{cases} -\frac{dp}{dt}(t) = p(t)F(t) - \frac{\partial l}{\partial x}(x(t), u(t), t) \\ p(T) = -g'(x(T)) \end{cases}$$

we can write:

$$\begin{aligned} \frac{\partial J}{\partial \theta_1} &= l(x(\theta_1), u(\theta_1^-), \theta_1) - l(x(\theta_1), u(\theta_1^+), \theta_1) + \\ &+ g'(x(T)) \frac{\partial}{\partial \theta_1} x(T) + \int_{\theta_1}^T \left(p(t)F(t) + \frac{dp}{dt}(t) \right) \frac{\partial}{\partial \theta_1} x(t) dt \end{aligned}$$

and integrating by parts we obtain:

$$\frac{\partial J}{\partial \theta_1} = l(x(\theta_1), u(\theta_1^-), \theta_1) - l(x(\theta_1), u(\theta_1^+), \theta_1) + p(\theta_1)G(\theta_1)(v_1 - v_2)(-1)^i$$

§4. CONTINUITY OF $\frac{\partial J}{\partial \theta_1}$

$$(1) \frac{\partial J}{\partial \theta_1} = l(x(\theta_1), u(\theta_1^-), \theta_1) - l(x(\theta_1), u(\theta_1^+), \theta_1) + p(\theta_1)G(\theta_1)(v_1 - v_2)(-1)^i .$$

We suppose that $F(\cdot), G(\cdot)$ are continuous and l, g are continuously differentiable. $\frac{\partial J}{\partial \theta_1}$ is a continuous function of $x(\theta_1), p(\theta_1)$ and θ_1 . The values $u(\theta_1^+)$ and $u(\theta_1^-)$ are constant (v_1 and v_2) , then we must prove only the continuity of $x(\theta_1), p(\theta_1)$ to obtain the continuity of $\frac{\partial J}{\partial \theta_1}$.

We can easily see that the transformation $\theta : \rightarrow u(\cdot)$ defined by

$$(2) \quad \begin{cases} u(t) = v_1 & 0 \leq t < \theta_1 \\ u(t) = v_1 + (v_2 - v_1) \left[\frac{1 - (-1)^i}{2} \right] & \theta_i \leq t < \theta_{i+1} \\ u(t) = v_1 + (v_2 - v_1) \left[\frac{1 - (-1)^n}{2} \right] & \theta_n \leq t \leq T \end{cases}$$

is continuous from $\Omega \rightarrow L_1(0, T)$.

$$(3) \quad x(t) = \Phi(t, 0)x_0 + \int_0^t \Phi(t, s)G(s)u(s)ds$$

and this formula defines a continuous transformation from $L_1(0, T) \rightarrow C(0, T; R^V)$ because, if $u_1(\cdot), u_2(\cdot)$ are two controls in $L_1(0, T)$ and $x_1(\cdot), x_2(\cdot)$ the system's evolution, it is:

$$(4) \quad \|x_1(t) - x_2(t)\| \leq M \int_0^T |u_1(s) - u_2(s)| ds = M \|u_1 - u_2\|_{L_1(0, T)}$$

where

$$(5) \quad M = \sup_{\substack{t \in [0, T] \\ s \in [0, t]}} \|\Phi(t, s)\| \cdot \sup_{s \in [0, T]} \|G(s)\|$$

and M is finite due to the continuity of G and Φ .

Obviously

$$(6) \quad x(\cdot) \rightarrow x(\theta_i) \text{ is continuous from } C(0, T; R^V) \rightarrow R^V .$$

Then,

$$(7) \quad \theta \rightarrow x(\theta_i) \text{ is continuous from } R^V \rightarrow R^V .$$

From the differential equation of $p(t)$ (p is a row vector)

$$(8) \quad - \frac{d}{dt} p(t) = p(t)F(t) - \frac{\partial 1}{\partial x}$$

we obtain the complete solution using the homogeneous solution and "variation of constants" method

$$(9) \quad p(t) = p(T)\Phi(T, t) + \int_T^t \frac{\partial 1}{\partial x}(x(s), u(s), s) \cdot \Phi(s, t) ds$$

$$(10) \quad p(T) = -g'(x(T))$$

thus, $-g'(x(T))\Phi(T, \theta_i)$ defines a continuous function from $\theta \in R^n \rightarrow R^V$, because g' is continuous, $x(T)$ is a continuous function of θ and Φ is absolutely continuous in its both arguments. The transformation $R^n \rightarrow L_1(0, T; R^V)$ given by:

$$\begin{array}{l} \theta \rightarrow u(\cdot) \rightarrow \frac{\partial 1}{\partial x}(x(t), u(t), t) \in L_1(0, T; R^V) \\ \quad \searrow \nearrow \\ \quad \quad \quad \rightarrow x(\cdot) \end{array}$$

is continuous from $R^n \rightarrow L_1(0, T; R^V)$.

To prove it, let θ, θ_ϵ be two set of switching points such that $\|\theta - \theta_\epsilon\| \rightarrow 0 \quad \epsilon \rightarrow 0$.

Then, $u(t) \rightarrow u_\epsilon(t) \quad \text{a.e.}$

and $\max_{t \in [0, T]} \|x(t) - x_\epsilon(t)\| \rightarrow 0$

now, $\frac{\partial 1}{\partial x}(x_\epsilon(s), u_\epsilon(s), s) \rightarrow \frac{\partial 1}{\partial x}(x(s), u(s), s) \quad \text{a.e.}$, and

$$\int_0^T \left\| \frac{\partial 1}{\partial x}(x_\epsilon(s), u_\epsilon(s), s) - \frac{\partial 1}{\partial x}(x(s), u(s), s) \right\| ds \rightarrow 0 \quad \text{as we can see applying}$$

Lebesgue's theorem.

The formula $\int_T^t \frac{\partial J}{\partial x}(x(s), u(s), s) \Phi(s, t) ds$ defines a continuous transformation from $L_1(0, T; R^V) \rightarrow C(0, T; R^V)$. Then, taking into account all these results we prove that $p(\theta_i)$ and also $\frac{\partial J}{\partial \theta_i}$ are continuous functions of θ .

§5. NUMERICAL SOLUTION OF THE PROBLEM BY THE APPLICATION OF THE PROJECTED GRADIENT METHOD.

We have seen that the problem of finding an optimal bang-bang policy with $n' \leq n$ switchings was reduced to the finite dimensional problem.

$$(1) \quad \min_{\Omega} J(\theta) \quad \Omega = \left\{ \theta \in R^n / 0 \leq \theta_1 \leq \theta_2 \leq \dots \leq \theta_n \leq T \right\}$$

Ω is a convex and compact subset where J is continuous, and this implicates the existence of a minimum.

We write the $n+1$ restrictions defining Ω in vector form.

$$f(\theta) \leq 0 \quad f \in R^{n+1} \quad \text{and} \quad f_1 = -\theta_1, f_2 = \theta_1 - \theta_2, \dots, f_n = \theta_{n-1} - \theta_n, f_{n+1} = \theta_n - T.$$

DEFINITION.

$$I_\epsilon(\theta) = \{i / f_i(\theta) + \epsilon \geq 0\}.$$

DEFINITION.

Given a set of integers, $I \subset \{1, \dots, n+1\}$, the projection of $y \in R^n$ on the subspace generated by the vectors ∇f_i , $i \in I$ is the vector $F_I \bar{\mu}$ that minimizes $\|y - F_I \mu\|^2$, where

$$F_I = [\nabla f_{i_1}, \dots, \nabla f_{i_m}] \quad , \quad I = \{i_1, \dots, i_m\} \quad \text{and} \quad \mu \in R^m.$$

It is easily shown that

$$\bar{\mu} = (F_I' F_I)^{-1} F_I' y \quad (F_I' \text{ is transpose of } F_I).$$

Then, the projection of y is:

$$P_I y = F_I (F_I' F_I)^{-1} F_I' y$$

and we can define the projection matrix

$$P_I = F_I (F_I' F_I)^{-1} F_I'.$$

In the same form we define the projection on the subspace orthogonal to ∇f_i , $i \in I$ and the corresponding matrix is:

$$P_I^\perp = I - P_I.$$

In the definition of P_I , we have supposed that the vectors ∇f_i , $i \in I$ are linearly independent, and then the matrix $F_I' F_I$ is invertible.

It is known (Kuhn-Tucker's theorem) that if $\bar{\theta}$ is a solution of the problem (1), then

$$\nabla J(\bar{\theta}) = F_{I_0}'(\bar{\theta}) \bar{\mu} \quad \bar{\mu}_1 \leq 0, \dots, \bar{\mu}_m \leq 0$$

and

$$\bar{\mu} = \left(F_{I_0}'(\bar{\theta}) F_{I_0}(\bar{\theta}) \right)^{-1} F_{I_0}'(\bar{\theta}) \nabla J(\bar{\theta})$$

We define a point as desirable if:

$$a) \quad \theta \in \Omega \quad b) \quad \nabla J(\theta) = F_{I_0}'(\theta) \mu(\theta) \quad \mu(\theta) \leq 0.$$

It is possible to apply the following algorithm: "Gradient Projected".

ALGORITHM.

- Step 0 : Select $\theta_0 \in \Omega$;
 $\varepsilon' > 0 / \forall \varepsilon > 0, \varepsilon \leq \varepsilon' \left\{ \nabla f_i(\theta) / i \in I_\varepsilon(\theta), \theta \in \Omega \right\}$ is a set of linearly independent vectors.
 Choose $\beta \in (0,1), \bar{\varepsilon} \in (0,\varepsilon'), \varepsilon'' \in (0,\bar{\varepsilon})$.
 Set $i = 0$.
- Step 1 : Set $\theta = \theta_i$.
- Step 2 : Set $\varepsilon_0 = \bar{\varepsilon}$ and $j = 0$.
- Step 3 : Compute $h_{\varepsilon_j} = P_{I_{\varepsilon_j}}^+(\theta) \nabla J(\theta)$.
- Step 4 : If $\|h_{\varepsilon_j}\| > \varepsilon_j$ $h(\theta) = -h_{\varepsilon_j}$ and go to step 12; else, go to 5.
- Step 5 : If $\varepsilon_j \leq \varepsilon''$, compute $h_0(\theta) = P_{I_0}^+(\theta) \nabla J(\theta)$ and
 $\mu_0(\theta) = \left(F_{I_0}^+(\theta) F_{I_0}(\theta) \right)^{-1} F_{I_0}^+(\theta) \nabla J(\theta)$ and go to 6; else, go to 7 .
- Step 6 : If $\mu_0(\theta) \leq 0$ and $\|h_0(\theta)\| = 0$ set $\theta_{i+1} = \theta$ and stop, else, go to 7.
- Step 7 : Compute $\mu_{\varepsilon_j}(\theta) = \left(F_{I_{\varepsilon_j}}^+(\theta) F_{I_{\varepsilon_j}}(\theta) \right)^{-1} F_{I_{\varepsilon_j}}^+(\theta) \nabla J(\theta)$.
- Step 8 : If $\mu_{\varepsilon_j}(\theta) \leq 0$, set $\varepsilon_{j+1} = \beta \varepsilon_j$, set $j = j+1$ and go to step 3 ; else, go to step 9 .
- Step 9 : Assuming that $I_{\varepsilon_j}(\theta) = \{k_1, \dots, k_{m'}\}$ and that $k_1 < k_2 < \dots < k_{m'}$, set $y_{\varepsilon_j}^\alpha(\theta) = \mu_{\varepsilon_j}^\alpha(\theta)$ for $\alpha = 1, 2, \dots, m'$ (where $\mu_{\varepsilon_j}^\alpha(\theta)$ is the α^{th} component of the vector $\mu_{\varepsilon_j}(\theta)$) .
- Step 10 : Find the smallest $k \in I_{\varepsilon_j}(\theta)$ such that the vector $\bar{h}_{\varepsilon_j}(\theta) = P_{I_{\varepsilon_j}}^+(\theta) - k \nabla J(\theta)$ satisfies the relation
 $\|\bar{h}_{\varepsilon_j}(\theta)\| = \max \left\{ \|P_{I_{\varepsilon_j}}^+(\theta) - 1 \nabla J(\theta)\| / 1 \in I_{\varepsilon_j}(\theta), y_{\varepsilon_j}^1(\theta) > 0 \right\}$
 and set $h(\theta) = -\bar{h}_{\varepsilon_j}(\theta)$.
- Step 11 : If $\|h(\theta)\| \leq \varepsilon_j$ set $\varepsilon_{j+1} = \beta \varepsilon_j$, set $j = j+1$, and go to step 3; else, go to step 12.
- Step 12 : Compute $\lambda(\theta) > 0$ to be the smallest scalar satisfying $J(\theta + \lambda(\theta)h(\theta)) = \min \{J(\theta + \lambda h(\theta)) / \lambda \geq 0, (\theta + \lambda h(\theta)) \in \Omega\}$.
- Step 13 : Set $\theta_{i+1} = \theta_i + \lambda(\theta)h(\theta)$, set $i = i+1$, and go to step 1 .

We have shown that J is continuously differentiable, then is valid the following theorem:

THEOREM. *The sequence θ_i given by the algorithm is finite and its last element is desirable or is infinite and each accumulation point of the sequence is desirable.*

(The proof of this Th. is, essentially, the same that we find in [1], pag. 195).

§6. NUMERICAL SOLUTION OF AN EXAMPLE: THE SHUT DOWN OF A NUCLEAR REACTOR.

The problem is the reduction of the power of a nuclear reactor in a fixed time. The functional to minimize is the xenon poisoning.

The model is ruled by the differential equations:

$$(1) \quad \begin{cases} \dot{I} = -aI + b\phi \\ \dot{x} = aI + c\phi - (d + e\phi)x \\ \dot{\phi} = U\phi. \end{cases}$$

I is the iodine concentration and x the xenon concentration.

ϕ is the flux of neutrons.

U is the control and it can only assume two values.

The control is applied in the interval $[0, T]$ in such a form that $\phi(T) = \phi_f$ (a fixed value). After that ($t > T$), the flux is held constant.

If we define $x_M = \max_{t > T} x(t)$, is possible to state the problem in the following form:

Find $U(t)$, $0 \leq t \leq T$, where $U(t)$ is a step function with n' switching ($n' \leq n$, n fixed), that takes only the values V_1, V_2 and such that the corresponding response of the system (1) satisfies $\phi(T) = \phi_f$ and gives the minimum value of $x_M(u(0, T))$.

This problem differs from the models studied in the nonlinearities of the equations (1) and in the fixed final condition.

In this case, it can be shown that $J(\theta) = x_M$ (where $\theta = (\theta_1, \dots, \theta_n)$ is the set of switching points) is a continuous function both with the derivatives of J , and then the theorem remains valid.

The final condition could be introduced in the functional through a penalization function. Another method is the following, we use the property that $\phi(T) = \phi_f$ implies that $\forall u(0, T) / \phi(T) = \phi_f$ $f^\circ(\theta) = \sum_{i=1}^{[(N+1)/2]} (\theta_{2i} - \theta_{2i-1}) = \text{constant}$

and consider this relation as an additional restriction. In the projected gradient algorithm, the matrix $F_{I_e}(\theta)$ is enlarged in the following form:

$$F_{I_e}(\theta) \rightarrow \tilde{F}_{I_e}(\theta) = \begin{bmatrix} \nabla f^\circ & \nabla f_{i_1} & \dots & \nabla f_{i_m} \end{bmatrix} \quad \{i_1, \dots, i_m\} = I_e(\theta)$$

and the points constructed by the algorithm satisfy $\phi(T) = \phi_f$, provided the initial point (θ_0) satisfies that condition.

FORMULAS OF J AND ∇J :

We have defined $J = x_M$; to compute it, we solve the equations:

$$\begin{cases} \dot{I} = -I + b\phi_f \\ \dot{x} = I + c\phi_f - (d + e\phi_f)x \end{cases} \quad \text{for } t > T$$

with initial conditions $x(T), I(T)$ and find the value $x_M = \max_{t \geq T} x(t)$. (We set $a=1$ making a change of variables).

If $\phi_f = 0$, it is:

$$\begin{cases} x_M = \exp(-t^*) \cdot \frac{I(T)}{d-1} + \exp(-dt^*) \left[x(T) - \frac{I(T)}{d-1} \right] & \text{if } \frac{x(T)}{I(T)} d < 1 \quad \text{and} \\ x_M = x(T) & \text{if } \frac{x(T)}{I(T)} d \geq 1 \quad \text{where} \\ t^* = \frac{1}{d-1} \ln \left[d + \frac{x(T)}{I(T)} d(1-d) \right] \end{cases} .$$

To compute $\frac{\partial J}{\partial \theta_1} = p_3(\theta_1) x_3(\theta_1) (v_2 - v_1) (-1)^i$, we integrate backwardly the adjoint equations:

$$\begin{cases} -\frac{dp_1}{dt} = -p_1 + p_2 \\ -\frac{dp_2}{dt} = -p_2(d + e\phi) \\ -\frac{dp_3}{dt} = p_1 + cp_2 - p_2 \cdot e \cdot x + p_3 U \end{cases}$$

with final conditions:

$$\begin{cases} p_1(T) = -\frac{\partial x_M}{\partial x(T)} \\ p_2(T) = -\frac{\partial x_M}{\partial I(T)} \\ p_3(T) = -\frac{\partial x_M}{\partial \phi(T)} \end{cases} .$$

The projected gradient algorithm (with the shown modification) was used to solve numerically the problem. The values of J and ∇J were computed integrating the differential equations of x, I, ϕ, p with a 4th order Runge - Kutta method.

NUMERICAL VALUES.

$$\begin{aligned} a &= 0.1 & \phi_f &= 0.674 \times 10^{-2} \\ b &= 1.0 & x_o &= 2.0 \\ c &= 1.0 & I_o &= 10.0 \\ d &= 0.05 & \phi_o &= 1.0 \\ e &= 0.95 & T &= 10.0 \\ v_1 &= -2.0 \\ v_2 &= 0.0 \end{aligned}$$

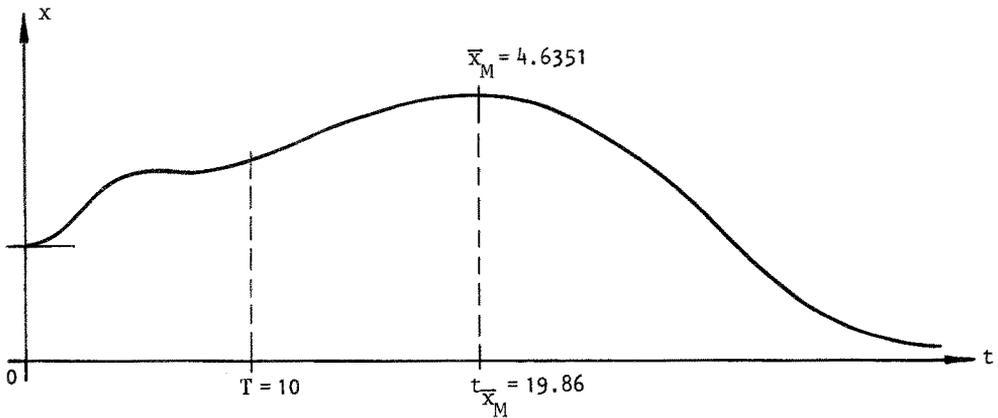
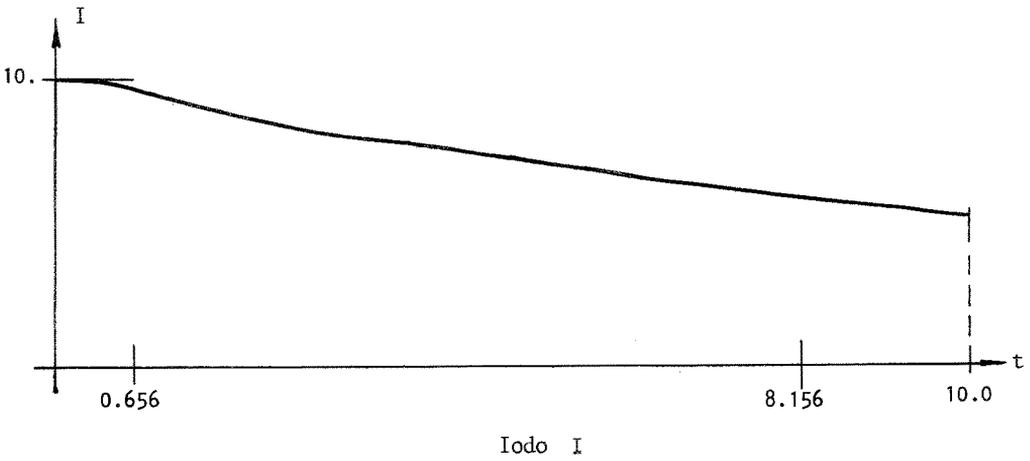
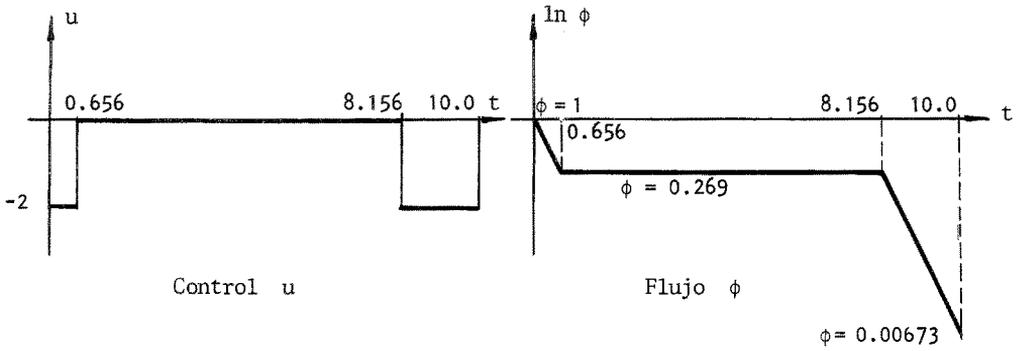
NUMERICAL RESULTS.

The optimal values obtained for $n=2$ are:

$$\begin{aligned} \bar{\theta}_1 &= 0.6561 \\ \bar{\theta}_2 &= 8.156 \\ \bar{x}_M &= 4.6351 \quad , \quad (\text{for } t_{\bar{x}_M} = 19.86) \end{aligned}$$

The following is a sample of the sequence produced by the algorithm, that shows rate of convergence.

θ_1	θ_2	x_M
1.000	8.500	4.82751
0.625	8.125	4.63729
0.659	8.154	4.63514
0.656	8.156	4.63513



Xenon: x

For $n \geq 2$ the sequences obtained were convergent to the optimal point for $n=2$, as is shown in the case $n=6$. Then, in the set of control with 6 or less switchings, the optimal one is a policy with two switchings.

θ_1	θ_2	θ_3	θ_4	θ_5	θ_6	X_M
0.250	1.000	1.250	2.000	2.250	8.250	4.7485
0.401	0.999	1.213	2.115	2.115	8.115	4.6634
0.488	0.993	1.195	2.115	2.115	8.191	4.6536
0.512	1.028	1.142	2.115	2.115	8.126	4.6470
0.568	1.026	1.132	2.115	2.115	8.172	4.6427
0.585	1.050	1.102	2.115	2.115	8.137	4.6396
0.645	1.059	1.085	2.115	2.115	8.171	4.6370
0.649	1.072	1.072	2.115	2.115	8.149	4.6352
0.656	1.072	1.072	2.115	2.115	8.156	4.6351

The switchings 2-3 and 4-5 are simultaneous and could be eliminated and the new policy is the optimal one for the problem with $n=2$.

57. FORM AND PROPERTIES OF THE SECOND DERIVATIVES OF J .

$$\frac{\partial J}{\partial \theta_1} = 1(x(\theta_1), u(\theta_1^-), \theta_1) - 1(x(\theta_1), u(\theta_1^+), \theta_1) + p(\theta_1)G(\theta_1)(v_1 - v_2))(-1)^i$$

a) $\frac{\partial^2 J}{\partial \theta_1^2}$

First, we find the formulas of $\frac{\partial x(\theta_1)}{\partial \theta_1}$, $\frac{\partial p(\theta_1)}{\partial \theta_1}$

$$x(\theta_1) = \Phi(\theta_1, 0)x_0 + \int_0^{\theta_1} \Phi(\theta_1, s)G(s)u(s)ds$$

then,

$$\frac{\partial x(\theta_1)}{\partial \theta_1} = F(\theta_1)x(\theta_1) + G(\theta_1)u(\theta_1^-)$$

Also,

$$p(\theta_1) = p(T)\Phi(T, \theta_1) - \int_{\theta_1}^T \frac{\partial 1}{\partial x}(x(s), u(s), s)\Phi(s, \theta_1)ds$$

with

$$p(T) = -\frac{\partial g}{\partial x}(x(T))$$

$$\begin{aligned} \frac{\partial p(\theta_1)}{\partial \theta_1} = & -p(T)\Phi(T, \theta_1)F(\theta_1) + \int_{\theta_1}^T \frac{\partial 1}{\partial x}(x(s), u(s), s)\Phi(s, \theta_1)F(\theta_1)ds + \\ & + \frac{\partial 1}{\partial x}(x(\theta_1), u(\theta_1^+), \theta_1) - \left(\frac{\partial}{\partial \theta_1}x(T)\right)' \frac{\partial^2 g}{\partial x^2}(x(T))\Phi(T, \theta_1) - \\ & - \int_{\theta_1}^T \left(\frac{\partial}{\partial \theta_1}x(s)\right)' \frac{\partial^2 1}{\partial x^2}(x(s), u(s), s)\Phi(s, \theta_1)ds \end{aligned}$$

but, we know that, for $t > \theta_1$

$$\frac{\partial}{\partial \theta_1}x(t) = \Phi(t, \theta_1)G(\theta_1)(v_2 - v_1)(-1)^i$$

then,

$$\begin{aligned} \frac{\partial p(\theta_1)}{\partial \theta_1} = & -p(\theta_1)F(\theta_1) + \frac{\partial 1}{\partial x}(x(\theta_1), u(\theta_1^+), \theta_1) - \\ & - (-1)^i(v_2 - v_1)G'(\theta_1)\Phi'(T, \theta_1) \frac{\partial^2 g}{\partial x^2}(x(T))\Phi(T, \theta_1) - \end{aligned}$$

$$- \int_{\theta_1}^T (-1)^i (v_2 - v_1) G'(\theta_1) \Phi'(s, \theta_1) \frac{\partial^2 1}{\partial x^2}(x(s), u(s), s) \cdot \Phi(s, \theta_1) ds .$$

Now, it is possible to compute $\frac{\partial^2 J}{\partial \theta_1^2}$

$$\begin{aligned} \frac{\partial^2 J}{\partial \theta_1^2} &= \frac{\partial}{\partial \theta_1} l(x(\theta_1), u(\theta_1^-), \theta_1) - \frac{\partial 1}{\partial \theta_1}(x(\theta_1), u(\theta_1^+), \theta_1) + \\ &+ \left[\frac{\partial}{\partial x} l(x(\theta_1), u(\theta_1^-), \theta_1) - \frac{\partial 1}{\partial x}(x(\theta_1), u(\theta_1^+), \theta_1) \right] \left[F(\theta_1)x(\theta_1) + G(\theta_1)u(\theta_1^-) \right] + \\ &+ p(\theta_1) \left(\frac{dG(\theta_1)}{d\theta_1} \right) (v_1 - v_2) (-1)^i + \left[-p(\theta_1)F(\theta_1) + \frac{\partial 1}{\partial x}(x(\theta_1), u(\theta_1^+), \theta_1) \right] . \\ &.G(\theta_1)(v_1 - v_2)(-1)^i + (-1)^i (v_1 - v_2) G'(\theta_1) \Phi'(T, \theta_1) \frac{\partial^2 g}{\partial x^2}(x(T)) \Phi(T, \theta_1) G(\theta_1) (v_1 - v_2) (-1)^i + \\ &+ (v_1 - v_2) (-1)^i G'(\theta_1) \int_{\theta_1}^T \Phi'(s, \theta_1) \frac{\partial^2 1}{\partial x^2}(x(s), u(s), s) \Phi(s, \theta_1) ds . G(\theta_1) (v_1 - v_2) (-1)^i . \end{aligned}$$

b) If $j > i$

$$\frac{\partial}{\partial \theta_j} x(\theta_1) = 0$$

then:

$$\frac{\partial}{\partial \theta_j} \left(\frac{\partial J}{\partial \theta_1} \right) = \frac{\partial}{\partial \theta_j} p(\theta_1) \cdot G(\theta_1) (v_1 - v_2) (-1)^i ,$$

and we must only know $\frac{\partial}{\partial \theta_j} p(\theta_1)$ to find the form of $\frac{\partial^2 J}{\partial \theta_j \partial \theta_1}$.

From the integral formula of p we obtain:

$$\begin{aligned} \frac{\partial}{\partial \theta_j} p(\theta_1) &= \left(\frac{\partial}{\partial \theta_j} x(T) \right)' \cdot \frac{\partial^2 g}{\partial x^2}(x(T)) \cdot \Phi(T, \theta_1) + \\ &+ \left[\frac{\partial 1}{\partial x}(x(\theta_j), u(\theta_j^+), \theta_j) - \frac{\partial 1}{\partial x}(x(\theta_j), u(\theta_j^-), \theta_j) \right] \Phi(\theta_j, \theta_1) - \\ &- \int_{\theta_j}^T (-1)^j (v_2 - v_1) G'(\theta_j) \Phi'(s, \theta_j) \frac{\partial^2 1}{\partial x^2}(x(s), u(s), s) \Phi(s, \theta_1) ds \end{aligned}$$

and then,

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \left(\frac{\partial J}{\partial \theta_1} \right) &= (-1)^j (v_2 - v_1) G'(\theta_j) \Phi(T, \theta_j) \frac{\partial^2 g}{\partial x^2}(x(T)) \cdot \Phi(T, \theta_1) G(\theta_1) (v_1 - v_2) (-1)^i + \\ &+ \left[\frac{\partial 1}{\partial x}(x(\theta_j), u(\theta_j^+), \theta_j) - \frac{\partial 1}{\partial x}(x(\theta_j), u(\theta_j^-), \theta_j) \right] \Phi(\theta_j, \theta_1) G(\theta_1) (v_1 - v_2) (-1)^i - \\ &- (-1)^j (v_2 - v_1) G'(\theta_j) \int_{\theta_j}^T \Phi'(s, \theta_j) \frac{\partial^2 1}{\partial x^2}(x(s), s) \Phi(s, \theta_1) ds \cdot G(\theta_1) (v_1 - v_2) (-1)^i \end{aligned}$$

c) If $i > j$

$$\frac{\partial}{\partial \theta_j} x(\theta_1) = \Phi(\theta_1, \theta_j) G(\theta_j) (v_2 - v_1) (-1)^j$$

and

$$\frac{\partial}{\partial \theta_j} p(\theta_1) = - \left(\frac{\partial}{\partial \theta_j} x(T) \right)' \frac{\partial^2 g}{\partial x^2}(x(T)) \Phi(T, \theta_1) - \int_{\theta_1}^T \left(\frac{\partial}{\partial \theta_j} x(s) \right)' \frac{\partial^2 1}{\partial x^2}(x(s), u(s), s) \Phi(s, \theta_1) ds$$

then,

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \left(\frac{\partial J}{\partial \theta_1} \right) &= \left[\frac{\partial 1}{\partial x}(x(\theta_1), u(\theta_1^-), \theta_1) - \frac{\partial 1}{\partial x}(x(\theta_1), u(\theta_1^+), \theta_1) \right] \Phi(\theta_1, \theta_j) . \\ &.G(\theta_j) (v_2 - v_1) (-1)^j + (-1)^j (v_2 - v_1) G'(\theta_j) \Phi'(T, \theta_j) \frac{\partial^2 g}{\partial x^2}(x(T)) \Phi(T, \theta_1) G(\theta_1) (v_1 - v_2) (-1)^i - \end{aligned}$$

$$- (-1)^j (v_2 - v_1) G'(\theta_j) \int_{\theta_i}^T \Phi'(s, \theta_j) \frac{\partial^2 \Phi}{\partial x^2}(x(s), u(s), s) \Phi(s, \theta_i) ds \cdot G(\theta_i) (v_1 - v_2) (-1)^i .$$

It can be proved, in the same form we have done for the first derivatives, that the second derivatives are continuous provided $g \in C^2$, $l \in C^2$, $G \in C^1$ and $F \in C$. This continuity is important to obtain superlinear convergence when it is applied the conjugate gradient method.

§8. SOLUTION OF THE PROBLEM USING A MIXED METHOD OF PENALIZATION AND CONJUGATE GRADIENTS.

The problem of minimum with restrictions:

$$(1) \quad \min_{\Omega} J(\theta) \quad ; \quad \Omega = \{ \theta / 0 \leq \theta_1 \leq \dots \leq \theta_n \leq T, \theta \in \mathbb{R}^n \}$$

is transformed into another that could be solved using the methods of optimization without restrictions. This is done applying penalty functions.

The new problem is:

$$\min_{\tilde{\Omega}} J_{\beta}(\theta) \quad ; \quad \tilde{\Omega} = \{ \theta \in \mathbb{R}^n / 0 < \theta_1 < \theta_2 < \dots < \theta_n < T \}$$

and

$$J_{\beta}(\theta) = J(\theta) + \beta \sum_{i=1}^{n+1} \varphi_i^{-1}(\theta)$$

$$\varphi_1 = -\theta_1 \quad ; \quad \varphi_i = -\theta_i + \theta_{i-1} \quad i = 2, \dots, n \quad ; \quad \varphi_{n+1} = \theta_n - T .$$

ALGORITHM.

Step 0: Choose $\theta_0 \in \tilde{\Omega}$; $\beta_0 > 0$; $\epsilon_0 > 0$ and set $i = 0$.

Step 1: Apply the conjugate gradient method to the minimization of J_{β} until it is obtained a point θ_{i+1} such that $\| \nabla J_{\beta}(\theta_{i+1}) \| < \epsilon_i$

Step 2: Let $\beta_{i+1} = \frac{\beta_i}{2}$, $\epsilon_{i+1} = \frac{\epsilon_i}{2}$, $i = i+1$ and go to 1.

REMARK.

The conjugate gradient method could be applied in $\tilde{\Omega}$, modifying the one dimensional search (along the conjugate directions) in such a form that the point are always chosen in $\tilde{\Omega}$.

It is known (Kuhn-Tucker's theorem) that a necessary condition for optimality of $\bar{\theta}$ in problem 1 is:

$$\nabla J(\bar{\theta}) + \sum_{i \in I(\bar{\theta})} \mu_i \nabla \varphi_i(\bar{\theta}) = 0 \quad \mu_i \geq 0$$

$$I(\theta) = \{ i / \varphi_i(\theta) = 0 \} .$$

DEFINITION.

A point θ is desirable if:

i) $\theta \in \Omega$

ii) $\nabla J(\theta) + \sum_{i \in I(\theta)} \mu_i \nabla \varphi_i(\theta) = 0$; $\mu_i \geq 0$ $i \in I(\theta)$

The algorithm has the property:

LEMA: If $J(\theta)$ is continuously differentiable, then the algorithm produces a sequence

of different points and all the accumulation points are desirable, or the sequence has a finite number of different points and the last (infinitely repeated) is desirable.

PROOF.

i) $\tilde{\Omega}$ is relatively compact, then $\{\theta_i\}$ has accumulation points.

Let be $\theta_{i_k} \rightarrow \tilde{\theta}$ $k = 1, 2, \dots$

$$\overline{\tilde{\Omega}} = \Omega, \text{ then } \tilde{\theta} \in \Omega.$$

It will be proved that $\tilde{\theta}$ satisfies the Kuhn-Tucker conditions.

From the step 1 of the algorithm it follows:

$$(1) \quad \nabla J(\theta_{i_k}) + \sum_{j=1}^{n+1} [\beta_{i_k} / \varphi_j^2(\theta_{i_k})] \nabla \varphi_j(\theta_{i_k}) \rightarrow 0.$$

It is easily seen that, $\forall \theta \in \Omega, I(\theta)$ has at most n elements and $\{\nabla \varphi_i(\theta), i \in I(\theta)\}$ is a set of linearly independent vectors.

Then,

$$\nabla J(\theta_{i_k}) + \sum_{j \in I(\tilde{\theta})} [\beta_{i_k} / \varphi_j^2(\theta_{i_k})] \nabla \varphi_j(\theta_{i_k}) \rightarrow 0.$$

We define: ψ_{i_k} (matrix $n \times m$)

$$\psi_{i_k} = [\nabla \varphi_{j_1}, \nabla \varphi_{j_2}, \dots, \nabla \varphi_{j_m}] \quad \{j_1, \dots, j_m\} = I(\tilde{\theta})$$

and then:

$$\begin{pmatrix} \beta_{i_k} / \varphi_{j_1}^2(\theta_{i_k}) \\ \vdots \\ \beta_{i_k} / \varphi_{j_m}^2(\theta_{i_k}) \end{pmatrix} - (\psi_{i_k}' \psi_{i_k})^{-1} \psi_{i_k}' \nabla J(\theta_{i_k}) \rightarrow 0$$

because, $\psi_{i_k} \rightarrow \psi = [\nabla \varphi_{j_1}(\tilde{\theta}), \dots, \nabla \varphi_{j_m}(\tilde{\theta})]$, this matrix is of maximum rank

and then $\psi_{i_k}' \psi_{i_k} \rightarrow \psi' \psi$ invertible matrix.

Then: $\beta_{i_k} / \varphi_{j_s}^2(\theta_{i_k}) \rightarrow \mu_{j_s} \geq 0 \quad j_s \in I(\tilde{\theta})$

and also, taking limits in (1):

$$\nabla J(\tilde{\theta}) + \sum_{j \in I(\tilde{\theta})} \mu_j \nabla \varphi_j(\tilde{\theta}) = 0.$$

Thus, $\tilde{\theta}$ is desirable.

ii) If $\exists N / \forall i \geq N, \theta_i = \theta_N$, from Step 1 it follows:

$$\left\| \nabla J(\theta_N) + \beta_i \nabla \left(\sum_{j=1}^{n+1} \varphi_j^{-1}(\theta_i) \right) \right\| \rightarrow 0 \quad i \rightarrow \infty$$

but $\beta_i \rightarrow 0$

then: $\nabla J(\theta_N) = 0$, $\theta_N \in \tilde{\Omega}$ and θ_N is desirable.

NUMERICAL RESULTS.

We have applied this method to the problem of the shutdown of the nuclear reactor. The results are the same obtained with the projected gradient method: the optimal policy in the set of step functions is a policy with two switchings.

The following table shows the convergence for $n=2$

θ_1	θ_2	J_β	β
1.000	8.500	6.62755	1.0000
0.913	8.413	5.49471	0.4000
0.792	8.292	4.98835	0.1600
0.720	8.220	4.77629	0.0640
0.684	8.184	4.69171	0.0256
0.667	8.167	4.65213	0.0076
0.659	8.159	4.64021	0.0023
0.658	8.158	4.63601	0.0004
0.6570	8.157	4.63530	0.00008
0.6564	8.1564	4.63516	0.000016

99. PROOF THAT THE SOLUTIONS OF THE PROBLEM WITH FIXED NUMBER OF SWITCHING ARE MINIMIZING SEQUENCE FOR THE PROBLEM WITH MEASURABLE CONTROLS.

The set of problems with step function controls (as were stated in §1) could be considered as a set of approximations to the problem:

$$\min_{\mathcal{U}_{ad}} J(u(.))$$

with

$$\mathcal{U}_{ad} = \left\{ u(.) \text{ measurable in } [0, T] / u(t) = v_1 \text{ or } u(t) = v_2 \text{ a.e.} \right\}$$

$$J(u(.)) = \int_0^T l(x(s), u(s), s) ds + g(x(T))$$

and $x(s)$ satisfies:

$$\begin{cases} \frac{dx}{dt}(t) = F(t)x(t) + G(t)u(t) \\ x(0) = x_0 \end{cases}$$

Under the assumptions that l, g are continuous and F, G are integrable, it can be proved that J is a continuous functional for $u \in \mathcal{U}_{ad}$ (with the $L_1(0, T)$ topology). If we denote with \bar{u}_n the optimal solution with at most n switchings, we shall prove that \bar{u}_n is a minimizing sequence for the new problem.

Let w_n be a minimizing sequence:

$$\lim_{n \rightarrow \infty} J(w_n) = \inf_{\mathcal{U}_{ad}} J(u) \quad w_n \in \mathcal{U}_{ad}$$

But, for the continuity of J , it is possible to find \tilde{w}_n (\tilde{w}_n a step function / $\tilde{w}_n \in \mathcal{U}_{ad}$) such that $|J(\tilde{w}_n) - J(w_n)| < \frac{1}{2^n}$

Let $k(n)$ be the number of switchings of \tilde{w}_n . By definition, $\forall k(n)$

$$J(\bar{u}_{k(n)}) \leq J(\tilde{w}_n)$$

Then, using the property that $J(\bar{u}_n)$ is non-increasing.

$$\lim_{m \rightarrow \infty} J(\bar{u}_m) \leq \lim_{n \rightarrow \infty} J(\tilde{w}_n) = \lim_{n \rightarrow \infty} J(w_n) = \inf_{\mathcal{U}_{ad}} J$$

but: $J(\bar{u}_m) \geq \inf_{\mathcal{U}_{ad}} J$

then $\lim_{m \rightarrow \infty} J(\bar{u}_m) = \inf_{u \in \mathcal{U}_{ad}} J(u)$.

§10. OTHERS RESULTS.

The bang-bang problem with restricted number of switchings could be analysed in global form (i.e., the initial state $x_0 \in \mathbb{R}^n$ or $x_0 \in \Omega$) and is reduced to a sequence of stopping-times problems. Also in the problem with measurable control ($u=0$ or $u=1$ a.e.) it is possible to prove existence theorems and to analyse the optimal cost function with the hamiltonian technique. These are the objects of forthcoming papers.

REFERENCES

1. E. POLAK. "Computational methods in optimization". Academic Press 1971.
2. E.B. LEE, L. MARKUS. "Foundations of Optimal Control Theory". Wiley. 1968.
3. L.S. PONTRYAGIN, V.G. BOLTYANSKII, R.V. GAMKRELIDZE, E.F. MISCHENKO. "Mathematical Theory of Optimal Processus". New York, Wiley 1962.-

Los originales de este trabajo fueron preparados en el Instituto de Matemática "Beppo Levi" por la Sra. H. I. Warecki de MUTY.