

State Space Synthesis of Discrete Linear Systems

VLADIMÍR STREJC

The target of this paper is to draw attention to some important relations between different procedures of state space synthesis of discrete linear systems optimizing the quadratic cost function. Common form of the synthesis concerns one-dimensional and multi-dimensional control problems solved by Second Method of Lyapunov, Maximum Principle of Pontryagin and by Bellman's Dynamic Programming. Emphasis is upon providing a comprehensive coverage that stresses general principles and integrates the well proved procedures into the overall picture.

1. INTRODUCTION

Without any doubt the quadratic cost function is one of the most useful mathematical gauges of the quality of automatic control. The procedures of synthesis using the quadratic cost function formerly elaborated for linear systems described in the s -plane by transfer functions represent today more or less classical approach developed for continuously and discontinuously acting one-dimensional and multi-dimensional systems, for analytical and random inputs and outputs and for command control and compensation of disturbances. The same is valid for the recent development of the synthesis in the state space concerning namely the continuously acting systems. The discrete version of this synthesis appeared only for some of the possible approaches and comparison of individual concepts is still lacking.

The aim of this article is to fill up the gaps and to compare the results obtained by different procedures, namely by the Second Methods of Lyapunov, Maximum Principle of Pontryagin and by Bellman's Dynamic Programming. All these methods are applied in this paper to discrete linear systems or to continuously acting linear systems controlled by digital computer provided that the system to be controlled is described by state space equations and that all state variables are measurable. The attention is drawn to the different types of Riccati equation and Euler equations in the case of the Maximum Principle of Pontryagin. The final calculation is assumed to be performed by solving the Riccati equation or by applying the calculated eigenvalues and eigenvectors of the fundamental matrix of the controlled system.

In order to have the final forms of the Riccati equation for the continuous version of systems and to have the possibility to compare them with the respective forms for the discrete version, the Second Method of Lyapunov was selected to demonstrate these relations.

2. GENERAL REMARKS

a) Assume a nonstationary linear controlled plant described in the state space by equations

$$(1) \quad \dot{\xi}(t) = A\xi(t) + B y(t),$$

$$(2) \quad x(t) = C\xi(t) + D y(t),$$

or in the common form

$$(3) \quad \dot{\xi}(t) = f_c[\xi(t), y(t), t],$$

$$(4) \quad x(t) = \varphi_c[\xi(t), y(t), t].$$

In the eqs. (1) through (4), $\xi(t)$ is a vector of state variables, $y(t)$ vector of controlling variables and $x(t)$ is a vector of output (controlled) variables.

Dimensions of the matrices are $A(n; n)$, $B(n; r)$, $C(p; n)$, $D(p; r)$ where n is the order of the controlled plant and $r \leq n$ and $p \leq n$. The elements of the matrices A , B , C , D are linear functions of the independent time variable t .

Let the rank of the matrix B be just r , i.e. its columns are linearly independent and let rank of the matrix C be just p , i.e. its rows are linearly independent.

The discrete version of the equations (1) and (2) has the form

$$(5) \quad \xi_{k+1} = F\xi_k + G y_k,$$

$$(6) \quad x_k = C\xi_k + D y_k,$$

or

$$(7) \quad \xi_{k+1} = f_d(\xi_k, y_k, k),$$

$$(8) \quad x_k = \varphi_d(\xi_k, y_k, k).$$

The dimensions of the matrices in eqs. (5) and (6) can be denoted in the same way as in the continuous version and the elements of these matrices are functions of the independent time variable k .

Solving eq. (1), we obtain

$$(9) \quad \xi(t) = F(t, t_0) \xi(0) + \int_{t_0}^t F(t, \tau) B(\tau) y(\tau) d\tau$$

where

$$(10) \quad F(t, t_0) = \exp \int_{t_0}^t A(\tau) d\tau$$

is the transition matrix of the time-varying system. At $t = t_0$, $F(t, t_0) = I$, where I is the identity matrix.

If the given plant is controlled by a digital computer, then the controlling variable is assumed to be constant between two successive intervals of sampling. Hence

$$(11) \quad y(\tau) = y_k \quad \text{for} \quad t_k \leq \tau \leq t_{k+1}$$

and

$$(12) \quad \xi_{k+1} = F(t_{k+1}, t_k) \xi_k + \int_{t_k}^{t_{k+1}} F(t_{k+1}, \tau) B(\tau) d\tau y_k = F_k \xi_k + G_k y_k$$

where

$$(13) \quad F_k = \exp \int_{t_k}^{t_{k+1}} A(\tau) d\tau = F[(k+1)T, kT],$$

$$(14) \quad G_k = \int_{t_k}^{t_{k+1}} F(t_{k+1}, \tau) B(\tau) d\tau.$$

For stationary processes, eqs. (13) and (14) reduce to

$$(15) \quad F = e^{AT},$$

$$(16) \quad G = \int_0^T e^{A\tau} B d\tau = A^{-1}(e^{AT} - I) B.$$

Let the matrices F_k or F respectively be nonsingular. The calculation of F from A is always unique if $\omega T \neq 2\pi\kappa$, $\kappa = 1, 2, \dots$, where T is the period of sampling.

b) The cost function is defined in the continuous case by the following form

$$(17) \quad J = \xi^T(t_1) P(t_1) \xi(t_1) + \int_{t_0}^{t_1} [\xi^T(t) Q(t) \xi(t) + y^T(t) R(t) y(t)] dt = \\ = \Theta[\xi(t_1), t_1] + \int_{t_0}^{t_1} \Phi[\xi(t), y(t), t] dt$$

and in the discrete version by

$$(18) \quad J = \xi_N^T P_N \xi_N + \sum_{k=0}^{N-1} (\xi_k^T Q_k \xi_k + y_k^T R_k y_k) = \\ = \Theta[\xi_N, N] + \sum_{k=0}^{N-1} \Phi[\xi_k, y_k, k]$$

where the relation in parenthesis in the second term of (18) represents the incremental cost for one stage of the discrete process.

The cost function can be applied to either a finite control interval $t_0 \leq t \leq t_1$ (or $0 \leq k \leq N$, $k = 0, 1, \dots, N$) or infinite interval $t_0 \leq t \leq \infty$ (or $0 \leq k \leq \infty$, $k = 0, 1, \dots, \infty$). This is in contrast to classical design procedures requiring the control interval to be infinity.

Matrices P, Q, R may be nonsymmetric but the same matrices being symmetric ones simplify essentially the resulting relations and the respective computational effort. For this reason only symmetric matrices P, Q, R are assumed in the next paragraphs. The matrix Q may be positive semidefinite. In some modifications of control problems the matrix R must be positive definite in order to ensure the existence of R^{-1} but it is not always the case and therefore attention will be paid to these particular solutions. The significance and the properties of the matrix P follow from the resulting relations.

Theorem 1. *If Q and R are symmetric and R positive definite then the functions Φ in (17) and (18) having the following more general form e.g.*

$$(19) \quad \Phi[\xi_k, y_k, k] = \xi_k^T Q_k \xi_k + 2\xi_k^T S_k y_k + y_k^T R_k y_k$$

can always be transformed into

$$(20) \quad \Phi[\hat{\xi}_k, \hat{y}_k, k] = \hat{\xi}_k^T \hat{Q}_k \hat{\xi}_k + \hat{y}_k^T \hat{R}_k \hat{y}_k$$

where

$$(21) \quad \begin{aligned} \hat{Q}_k &= Q_k - S_k R_k^{-1} S_k^T, \\ \hat{R}_k &= R_k, \\ \hat{\xi}_k &= \xi_k, \\ \hat{y}_k &= y_k + R_k^{-1} S_k^T \xi_k. \end{aligned}$$

Applying the transformation (21), the state equation changes into

$$(22) \quad \hat{\xi}_{k+1} = \hat{F}_k \hat{\xi}_k + \hat{G}_k \hat{y}_k$$

where

$$(23) \quad \begin{aligned} \hat{F}_k &= F_k - G_k R_k^{-1} S_k^T, \\ \hat{G}_k &= G_k. \end{aligned}$$

Proof. Starting with relation (19) and adding and subtracting the term $\xi_k^T S_k R_k^{-1} S_k^T \xi_k$, then, after some rearrangements, it is possible to prove the resulting relation (20).

c) Let the control law be

$$(24) \quad y(t) = M(t) \xi(t)$$

for the continuous version and

$$(25) \quad y_k = M_k \xi_k$$

for the discrete version.

d) The control law is optimal according to the quadratic cost function, if for any initial state $\xi(0) \in X_n$ of the controlled plant it holds that

1. The functional (17) or (18) respectively reaches its minimal value.
2. The control loop is stable.

Remark. For the discrete version controlled plants with shifted output are admitted.

e) The controlled plant described by eqs. (5) and (6) has a shifted output by m periods of sampling if for arbitrary initial state $\xi(0) \in X_n$ the output ξ_k , $k \geq m$, does not depend on inputs $y_k, y_{k-1}, \dots, y_{k-m+1}$ but depends on y_{k-m} and eventually on next past inputs.

In order to have a mathematical rule for the calculation of m , it is possible to eliminate successively the state vectors from (5) and (6) in the following way

$$(26) \quad \begin{aligned} \xi_1 &= F\xi_0 + Gy_0, \\ \xi_2 &= F^2\xi_0 + FGy_0 + Gy_1, \\ &\vdots \\ \xi_k &= F^k\xi_0 + F^{k-1}Gy_0 + \dots + FGy_{k-2} + Gy_{k-1}. \end{aligned}$$

Hence eq. (6) with (26) yields

$$(27) \quad x_i = CF^k\xi_0 + CF^{k-1}Gy_0 + \dots + CGy_{k-1} + Dy_k.$$

Theorem 2. The controlled plant has the output shifted by m periods of sampling if $\Gamma_i = 0$ for $i = 0, 1, \dots, m-1$ and $\Gamma_m \neq 0$, where $\Gamma_0 = D$ and $\Gamma_i = CF^{i-1}G$.

3. CONTINUOUS OPTIMAL CONTROL VIA SECOND METHOD OF LYAPUNOV

3.1 General procedure

In this section we consider the problem of calculating the control vector $y(t)$ so as to minimize the cost function of the system to be transferred from the initial state

$\xi_0 \neq 0$ at $t = t_0$ as close as possible to the desired terminal state, the origin of the state space, by applying the control vector $y(t)$ to the plant. The problem is to be solved for finite control interval having fixed beginning and terminal times. Since the problem is considered to be linear one no inequality constraints will be applied.

The second method of Lyapunov attempts to give information on the stability of equilibrium state of linear and nonlinear systems without any knowledge of their solutions and consists of determination of a fictitious function called Lyapunov function $V[\xi(t), t]$ the sign of which and the sign of its time derivative $\dot{V}[\xi(t), t]$ enables to check the stability of the equilibrium state under consideration. Without going into the details described in the technical literature it may be pointed out for the purpose of this article that the Lyapunov function $V[\xi(t), t]$ is a scalar positive definite function and it is continuous together with its first partial derivatives with respect to its arguments in region Ω about the origin and has a time derivative which is negative definite (or semidefinite). Notice that $\dot{V}[\xi(t), t]$ is actually the total derivative of $V[\xi(t), t]$ with respect to t and $\dot{V}[\xi(t), t] < 0$ implies that $V[\xi(t), t]$ is a decreasing function of t .

Theorem 3. *If the system is defined by equation*

$$\dot{\xi}(t) = f_c[\xi(t), t]$$

where $f_c(0, t) = 0$ for all t and if there exists scalar function $V[\xi(t), t]$, with continuous first partial derivatives, satisfying the following conditions

- a) $V[\xi(t), t] \geq \mathcal{V}[\xi(t)] > 0$ for all $\xi(t) \neq 0$ in Ω and all t ,
 $V[0, t] = 0$ for all t ,
- b) $\dot{V}[\xi(t), t] \leq \mathcal{W}[\xi(t)] < 0$ for all $\xi(t) \neq 0$ in Ω and all t ,
 $\dot{V}[0, t] = 0$ for all t

then the system is uniformly asymptotically stable in Ω .

The Lyapunov function is not unique for a given system and therefore the second method of Lyapunov can be used not only for stability considerations but for more general problems of synthesis.

For the control problem under consideration, let the Lyapunov function be

$$(28) \quad V[\xi(t), t] = \xi^T(t) P(t) \xi(t)$$

where P is a positive definite matrix.

Then

$$(29) \quad \int_{t_0}^{t_1} \frac{\partial}{\partial t} \xi^T(t) P(t) \xi(t) dt = \xi^T(t) P(t) \xi(t) \Big|_{t_0}^{t_1}$$

or

$$(30) \quad \int_{t_0}^{t_1} [\xi^T(t) P(t) \xi(t) + \xi^T(t) P(t) \dot{\xi}(t) + \xi^T(t) \dot{P}(t) \xi(t)] dt = \\ = \xi^T(t_1) P(t_1) \xi(t_1) - \xi^T(t_0) P(t_0) \xi(t_0).$$

Substituting $\dot{\xi}(t)$ from (1) into the last equation, it is possible to write

$$(31) \quad \int_{t_0}^{t_1} [\xi^T(t) A^T(t) P(t) \xi(t) + \xi^T(t) P(t) A(t) \xi(t) + \\ + y^T(t) B^T(t) P(t) \xi(t) + \xi^T(t) P(t) B(t) y(t) + \xi^T(t) \dot{P}(t) \xi(t)] dt - \\ - \xi^T(t_1) P(t_1) \xi(t_1) + \xi^T(t_0) P(t_0) \xi(t_0) = 0.$$

In order to ensure the asymptotic stability of the solution the time derivative of the Lyapunov function must be negative definite. As mentioned above the Lyapunov function is not unique and therefore it can be selected properly in such a way to create a relation with the cost function and to establish the necessary condition for the optimum control. The desired relation is obviously

$$(32) \quad \frac{\partial}{\partial t} \xi^T(t) P(t) \xi(t) = - [\xi^T(t) Q(t) \xi(t) + y^T(t) R(t) y(t)].$$

Evidently, if Q and R are symmetric matrices then the matrix P is symmetric matrix as well.

We wish to minimize the cost function (17). To solve the problem it is possible to start by adding relation (31), having zero value, to the cost function (17). We obtain

$$(33) \quad J = \xi^T(t_1) P(t_1) \xi(t_1) + \int_{t_0}^{t_1} [\xi^T(t) Q(t) \xi(t) + y^T(t) R(t) y(t) + \\ + \xi^T(t) A^T(t) P(t) \xi(t) + \xi^T(t) P(t) A(t) \xi(t) + \\ + y^T(t) B^T(t) P(t) \xi(t) + \xi^T(t) P(t) B(t) y(t) + \\ + \xi^T(t) \dot{P}(t) \xi(t)] dt - \xi^T(t_1) P(t_1) \xi(t_1) + \xi^T(t_0) P(t_0) \xi(t_0).$$

The minimum of (33) requires that

$$(34) \quad \frac{\partial J}{\partial y(t)} = 0$$

and $\partial^2 J / \partial y^2(t)$ to be a positive definite matrix. Condition (34) yields

$$(35) \quad R(t) y(t) + B^T(t) P(t) \xi(t) = 0$$

90 which defines the optimum control vector

$$(36) \quad y(t) = -R^{-1}(t) B^T(t) P(t) \xi(t)$$

with

$$(37) \quad M(t) = -R^{-1}(t) B^T(t) P(t)$$

being the feedback transition matrix.

Since

$$(38) \quad \frac{\partial^2 J}{\partial y^2(t)} = R(t)$$

is positive definite matrix, the cost function reaches with the optimum control vector (36) its minimum.

In order to find the final value of the cost function we can rewrite relation (32) into

$$(39) \quad \begin{aligned} & \int_{t_0}^{t_1} [\xi^T(t) P(t) \xi(t) + \xi^T(t) P(t) \dot{\xi}(t) + \xi^T(t) \dot{P}(t) \xi(t)] dt = \\ & = - \int_{t_0}^{t_1} [\xi^T(t) Q(t) \xi(t) + y^T(t) R(t) y(t)] dt \end{aligned}$$

and express the term

$$(40) \quad \begin{aligned} & \int_{t_0}^{t_1} \xi^T(t) \dot{P}(t) \xi(t) dt = - \int_{t_0}^{t_1} \{ \xi^T(t) [A^T(t) P(t) + P(t) A(t)] \xi(t) + \\ & + y^T(t) B^T(t) P(t) \xi(t) + \xi^T(t) P(t) B(t) y(t) + \xi^T(t) Q(t) \xi(t) + y^T(t) R(t) y(t) \} dt. \end{aligned}$$

Finally, substituting relation (40) into eq. (33) for the last term in the integrand, the integrand in (33) vanishes and the minimum value of the cost function is

$$(41) \quad J = \xi^T(t_0) P(t_0) \xi(t_0).$$

Hence, for the closed control loop holds

$$(42) \quad \dot{\xi}(t) = [A(t) - B(t) R^{-1}(t) B^T(t) P(t)] \xi(t).$$

If F is the transition matrix of this equation then the control vector is given by

$$(43) \quad y(t) = -R^{-1}(t) B^T(t) P(t) F(t, t_0) \xi(t_0).$$

From eqs. (36) it is obvious that the matrix $R(t)$ must be positive definite.

Unfortunately, the matrix $P(t)$ in the resulting relations is not yet known and must be calculated first. Substituting from (36) the control vector $y(t)$ into eq. (40) then

this equation holds for arbitrary $\xi(t)$ only if

$$(44) \quad \dot{P}(t) + A^T(t)P(t) + P(t)A(t) - P(t)B(t)R^{-1}(t)B^T(t)P(t) + Q(t) = 0.$$

This is the general form of the matrix differential Riccati equation the solution of which yields the desired matrix $P(t)$.

The final form of the control vector (36) requires the matrix R to be positive definite then the inverse R^{-1} must exist. This mathematical condition corresponds to the physical requirements that all controlling variables of the vector y in the case of continuous linear system must be constrained to obtain a physically realizable process. On the other hand if some of the controlling variables would not be constrained it would attain a Dirac impuls shape at the beginning of the process which is not realizable.

3.2 Simplified modifications

a) For $t_0 = 0$, t_1 approaching to infinity and for

$$(45) \quad \lim_{t_1 \rightarrow \infty} \xi^T(t_1)P(t_1)\xi(t_1) = 0$$

the cost function (17) takes the form

$$(46) \quad J = \int_0^\infty [\xi^T(t)Q(t)\xi(t) + y^T(t)R(t)y(t)] dt.$$

The problem described in this way i.e. the problem of minimization of system terminal errors, is usually called a noise-free optimal regulator problem.

Owing to (45) and (46) some of the relations of the preceding section are slightly simplified but final results remain unchanged.

b) If the performance index is not time weighted and if the controlled plant is time invariant then all matrices in the state eqs. (1) and (2) are constant matrices and particularly the matrix of the Lyapunov function $P(t) = P = \text{const.}$ Hence $\dot{P}(t) = 0$.

Clearly, for the noise-free regulator problem the Riccati eq. (44) reduces to algebraic quadratic matrix equation

$$(47) \quad A^T P + P A - P B R^{-1} B^T P + Q = 0.$$

Applying the introduced simplifications there is no difficulty to modify all other relations from section 3.1 describing the general case.

c) Considering only a homogeneous state equation of the time invariant system

$$(48) \quad \dot{\xi}(t) = A \xi(t)$$

92 and a constant weighting matrix Q in the cost function

$$(49) \quad I = \int_0^{\infty} \xi^T(t) Q \xi(t) dt$$

then with $\lim_{t \rightarrow \infty} \xi^T(t) P \xi(t) = 0$ eqs. (32) and (39) yield

$$(50) \quad \int_0^{\infty} \frac{\partial}{\partial t} \xi^T(t) P \xi(t) dt = -\xi^T(0) P \xi(0),$$

$$(51) \quad \int_0^{\infty} \xi^T(t) (A^T P + P A) \xi(t) dt = -\int_0^{\infty} \xi^T(t) Q \xi(t) dt.$$

Hence the Riccati equation is now

$$(52) \quad A^T P + P A + Q = 0.$$

It is a linear matrix equation having an explicit solution which may be used as the first estimate of the more general forms of Riccati equation.

Eq. (33) simplified for this particular case, gives the final value of the performance index

$$(53) \quad J = \xi^T(0) P \xi(0).$$

Clearly, there is no controlling action. Solution of eq. (48)

$$(54) \quad \xi(t) = e^{At} \xi(0)$$

corresponds to eigenoscillations of the plant itself when no input signal is acting on the system under consideration and the cost function (49) can be considered as a quality performance index of these oscillations.

4. DISCRETE OPTIMAL CONTROL PROBLEMS VIA SECOND METHOD OF LYAPUNOV

4.1 General procedure

This section presents the same problem as in section 3.1 however for discrete time systems. First introduce the discrete version of the second method of Lyapunov.

Theorem 4. If the system is defined by equation $\xi_{k+1} = f_d(\xi_k, k)$ where $f_d(0, k) = 0$ and if there exists a scalar function $V(\xi, k)$ continuous in ξ , such that

$$a) \quad V(\xi, k) \geq \mathcal{V}(\xi) > 0 \quad \text{for } \xi \neq 0 \quad \text{and all } k,$$

$$V(0, k) = 0 \quad \text{for all } k,$$

$$b) \quad \Delta V(\xi, k) \leq \mathcal{W}(\xi) < 0 \quad \text{for } \xi \neq 0 \quad \text{and all } k$$

where

$$\Delta V(\xi_k, k) = V(\xi_{k+1}, k) - V(\xi_k, k),$$

$$\text{c)} \quad V(\xi, k) \rightarrow \infty \quad \text{as} \quad \|\xi\| \rightarrow \infty$$

then the equilibrium state $\xi = 0$ is asymptotically stable in the large and $V(\xi, k)$ is a Lyapunov function.

For the discrete system the Lyapunov function may have the form

$$(55) \quad V(\xi, k) = \xi_k^T P_k \xi_k$$

where P is a positive definite matrix. Then

$$(56) \quad \Delta V(\xi, k) = \Delta(\xi_k^T P_k \xi_k) = \xi_{k+1}^T P_{k+1} \xi_{k+1} - \xi_k^T P_k \xi_k.$$

The difference of the Lyapunov function must be negative definite. Combining this condition with the summand of the cost function (18), the matrix P_k is essentially determined. We have

$$(57) \quad \xi_{k+1}^T P_{k+1} \xi_{k+1} - \xi_k^T P_k \xi_k = -(\xi_k^T Q_k \xi_k + y_k^T R_k y_k).$$

In order to minimize the cost function, let us calculate the sum of (56) first

$$(58) \quad \sum_{k=0}^{N-1} \Delta(\xi_k^T P_k \xi_k) = \sum_{k=0}^{N-1} (\xi_{k+1}^T P_{k+1} \xi_{k+1} - \xi_k^T P_k \xi_k),$$

$$\sum_{k=0}^{N-1} \Delta(\xi_k^T P_k \xi_k) = \xi_N^T P_N \xi_N - \xi_0^T P_0 \xi_0$$

or if ξ_{k+1} in (58) is expressed according to the state equation, then

$$(59) \quad \sum_{k=0}^{N-1} (\xi_{k+1}^T P_{k+1} \xi_{k+1} - \xi_k^T P_k \xi_k) =$$

$$= \sum_{k=0}^{N-1} (\xi_k^T F_k^T P_{k+1} F_k \xi_k + y_k^T G_k^T P_{k+1} F_k \xi_k + \xi_k^T F_k^T P_{k+1} G_k y_k +$$

$$+ y_k^T G_k^T P_{k+1} G_k y_k - \xi_k^T P_k \xi_k).$$

Comparing the right hand sides of eqs. (58) and (59) evidently the following equality holds

$$(60) \quad \sum_{k=0}^{N-1} (\xi_k^T F_k^T P_{k+1} F_k \xi_k + y_k^T G_k^T P_{k+1} F_k \xi_k + \xi_k^T F_k^T P_{k+1} G_k y_k +$$

$$+ y_k^T G_k^T P_{k+1} G_k y_k - \xi_k^T P_k \xi_k) - \xi_N^T P_N \xi_N + \xi_0^T P_0 \xi_0 = 0.$$

94 Now, the cost function (18) extended by (60) gives

$$(61) \quad J = \xi_N^T P_N \xi_N + \sum_{k=0}^{N-1} (\xi_k^T Q_k \xi_k + y_k^T R_k y_k) + \sum_{k=0}^{N-1} (\xi_k^T F_{k+1}^T P_{k+1} F_k \xi_k + y_k^T G_k^T P_{k+1} F_k \xi_k + \xi_k^T F_{k+1}^T P_{k+1} G_k y_k + y_k^T G_k^T P_{k+1} G_k y_k - \xi_k^T P_k \xi_k) - \xi_N^T P_N \xi_N + \xi_0^T P_0 \xi_0.$$

The optimum control vector must satisfy following condition

$$(62) \quad \frac{\partial J}{\partial y_k} = 0$$

provided that the second partial derivative i.e. $\partial^2 J / \partial y_k^2$ is a positive definite matrix. Condition (62) applied on relation (61) gives

$$(63) \quad (R_k + G_k^T P_{k+1} G_k) y_k + G_k^T P_{k+1} F_k \xi_k = 0.$$

Hence the optimum control vector is

$$(64) \quad y_k = -(G_k^T P_{k+1} G_k + R_k)^{-1} G_k^T P_{k+1} F_k \xi_k$$

where

$$(65) \quad M_k = -(G_k^T P_{k+1} G_k + R_k)^{-1} G_k^T P_{k+1} F_k$$

represents the feedback transition matrix.

Besides, $\partial^2 J / \partial y_k^2$ calculated by means of (63) yields

$$(66) \quad \frac{\partial^2 J}{\partial y_k^2} = R_k + G_k^T P_{k+1} G_k$$

which in accordance with the excepted assumptions is a positive definite matrix and consequently the control vector (64) ensures the minimum of the cost function.

The final value of the cost function may be calculated in the following way. Substituting into (57) for ξ_{k+1} from the state equation (5) and using for y_k the relation (25) of the control law, eq. (57) takes the form

$$(67) \quad \xi_k^T (F_k^T P_{k+1} F_k + M_k^T G_k^T P_{k+1} F_k + F_k^T P_{k+1} G_k M_k + M_k^T G_k^T P_{k+1} G_k M_k - P_k) \xi_k = -[\xi_k^T (Q_k + M_k^T R_k M_k) \xi_k].$$

Since eq. (67) must be satisfied for arbitrary ξ_k , it holds that

$$(68) \quad P_k = F_k^T P_{k+1} F_k + M_k^T G_k^T P_{k+1} F_k + F_k^T P_{k+1} G_k M_k + M_k^T G_k^T P_{k+1} G_k M_k + Q_k + M_k^T R_k M_k.$$

Now inserting into (61) for P_k relation (68) and for y_k the control law (25), then the minimum value of the cost function is

$$(69) \quad J = \xi_0^T P_0 \xi_0.$$

The difference equation of the closed control loop is

$$(70) \quad \xi_{k+1} = [F_k - G_k(G_k^T P_{k+1} G_k + R_k)^{-1} G_k^T P_{k+1} F_k] \xi_k.$$

Denoting the transition matrix of eq. (70) according to (13), the control vector has the final form

$$(71) \quad y_k = - (G_k^T P_{k+1} G_k + R_k)^{-1} G_k^T P_{k+1} F_{k+1,0} \xi_0.$$

Notice that in contrast to the continuous version of the problem under consideration it is not necessary for the matrix R_k to be positive definite. This mathematical result corresponds again to the physical reality then for a discrete linear system the controlling variables in vector y need not be constrained in order to achieve a physically realizable process. Hence the constraint applied on control vector y may concern some controlling variables only. Consequently, the weighting matrix R may be positive semidefinite. Eq. (68) with M_k according to (65) represents the matrix difference "Riccati* equation" of the form

$$(72) \quad P_k = F_k^T P_{k+1} F_k - F_k^T P_{k+1} G_k (G_k^T P_{k+1} G_k + R_k)^{-1} G_k^T P_{k+1} F_k + Q_k$$

the solution of which gives the matrix P_k needed for the numerical calculation of the control vector y_{k-1} . Notice P_{k+1} must be known if P_k is to be calculated. Hence, the matrices P_k can be determined starting with the endpoint of the process only.

4.2 Simplified modifications

a) For N approaching to infinity and for

$$(73) \quad \lim_{N \rightarrow \infty} \xi_k^T P_N \xi_N = 0$$

the cost function (18) changes into

$$(74) \quad J = \sum_{k=0}^{\infty} (\xi_k^T Q_k \xi_k + y_k^T R_k y_k)$$

and the problem solved according to this criterion corresponds again to a noise-free optimal regulator problem.

* Eq. (72) does not correspond exactly to the discrete matrix form of the Riccati equation but we use this denotation in order to stress the relation to the continuous version of this equation and to accentuate its significance. The same remark holds for eq. (102), (129), (146). On the other hand eq. (110) has a form very close to the usual one.

With simplifications (73) and (74) the final results described in paragraph 4.1 in a common way remain unchanged except that the matrix $P_k = P_{k+1} = P$ is a constant matrix.

b) For time invariant controlled plant and for the cost function not time weighted, it is possible to omit all indexes k and $k + 1$ at matrices in all relations of the section 4.1 and the resulting relationships remain valid. Notice, that in this case the Riccati equation simplifies to algebraic nonlinear matrix equation.

c) Applying simplifications formulated in a) and b) and in addition to it the weighting matrix $R = 0$, we have the cost function

$$(75) \quad J = \sum_{k=0}^{\infty} \xi_k^T Q \xi_k,$$

the control vector

$$(76) \quad y_k = - (G^T P G)^{-1} G^T P F \xi_k,$$

the closed control loop equation

$$(77) \quad \xi_{k+1} = [F - G(G^T P G)^{-1} G^T P F] \xi_k,$$

and the Riccati equation of the form

$$(78) \quad P = F^T P F - F^T P G (G^T P G + R)^{-1} G^T P F + Q.$$

d) Considering a homogeneous equation of a system and all other simplifications indicated in a) through c), then eq. (67) takes the form

$$(79) \quad \xi_k^T (F^T P F - P) \xi_k = - \xi_k^T Q \xi_k.$$

Hence the Riccati equation is a linear algebraic matrix equation of the form

$$(80) \quad F^T P F - P + Q = 0$$

and the minimum value of the cost function is given again by eq. (69). Similar remarks might be expressed here concerning this particular case as they were stated at the end of the section 3.2.

5. DISCRETE MAXIMUM PRINCIPLE

5.1. General procedure

One of the most useful techniques in modern control theory is that branch of mathematics known as the calculus of variations. There are two different procedures frequently applied for the synthesis of general control problems, the Euler-Lagrange technique and the maximum principle of Pontryagin, both differing in the

mathematical background. We shall draw our attention especially to the discrete version of the maximum principle and indicate common relations. We shall not apply for the purpose of this article a quite general formulation of the problem, which might be solved by maximum principle. On the other hand it is possible to derive the fundamental relations in such a way to be valid for linear and nonlinear systems as well.

Hence let the dynamic system be discrete and nonlinear one with the state vector ξ_k and the input vector y_k . The system is described by eqs (7) and (8) and the process terminates at stage N . The problem is to find y_k such as to minimize the cost function (18) subject to the constraint (7). Taking this constraint into account, the cost function (18) can be written in the general form

$$(81) \quad J_C = \Theta(\xi_k, N) \Big|_0^N + \sum_{k=0}^{N-1} \Phi(\xi_k, y_k, k) - \lambda_{k+1}^T [\xi_{k+1} - f_d(\xi_k, y_k, k)]$$

where λ is a vector of Lagrange multipliers.

The Hamiltonian is defined by

$$(82) \quad H(\xi_k, y_k, \lambda_{k+1}, k) = H_k = \Phi(\xi_k, y_k, k) - \lambda_{k+1}^T f_d(\xi_k, y_k, k).$$

With (82) the cost function then becomes

$$(83) \quad J_C = \Theta(\xi_k, N) \Big|_0^N + \sum_{k=0}^{N-1} (H_k - \lambda_{k+1}^T \xi_{k+1}).$$

Now to obtain the minimum of (83) with respect to ξ_k and y_k , the method of perturbations of the calculus of variations may be used. Hence, let us introduce for the state and input vectors following relations

$$(84) \quad \begin{aligned} \xi_k &= \hat{\xi}_k + \varepsilon \delta_k, \\ y_k &= \hat{y}_k + \varepsilon \eta_k, \end{aligned}$$

where the perturbations δ_k and η_k are mutually independent and their values at different stages are independent too.

With relations (84), the cost function (83) takes the form

$$(85) \quad \begin{aligned} J_C &= \Theta(\hat{\xi}_N + \varepsilon \delta_N, N) - \Theta(\hat{\xi}_0 + \varepsilon \delta_0, k_0) + \\ &+ \sum_{k=0}^{N-1} [H(\hat{\xi}_k + \varepsilon \delta_k, \hat{y}_k + \varepsilon \eta_k, k) - \lambda_{k+1}^T (\hat{\xi}_{k+1} + \varepsilon \delta_{k+1})] \end{aligned}$$

and the minimum of J_C requires that

$$(86) \quad \lim_{\varepsilon \rightarrow 0} \frac{\partial J_C}{\partial \varepsilon} = 0 \quad \text{and} \quad \lim_{\varepsilon \rightarrow 0} \frac{\partial^2 J_C}{\partial \varepsilon^2} > 0.$$

98 The first condition yields

$$(87) \quad \left(\frac{\partial \Theta_N}{\partial \xi_N} \right)^T \delta_N - \left(\frac{\partial \Theta_0}{\partial \xi_0} \right)^T \delta_0 + \sum_{k=0}^{N-1} \left[\left(\frac{\partial H_k}{\partial \xi_k} \right)^T \delta_k + \left(\frac{\partial H_k}{\partial \eta_k} \right)^T \eta_k - \lambda_{k+1}^T \delta_{k+1} \right] = 0.$$

The last term in eq. (87) can be rewritten as follows

$$(88) \quad \sum_{k=0}^{N-1} \lambda_{k+1}^T \delta_{k+1} = \sum_{k=1}^N \lambda_k^T \delta_k = \sum_{k=0}^{N-1} \lambda_k^T \delta_k + \lambda_N^T \delta_N - \lambda_0^T \delta_0.$$

Consequently, if again ξ_k and y_k is applied instead of ξ_k and η_k respectively, condition (87) with relation (88) may be expressed as

$$(89) \quad \left[\left(\frac{\partial \Theta_N}{\partial \xi_N} \right)^T - \lambda_N^T \right] \delta_N - \left[\left(\frac{\partial \Theta_0}{\partial \xi_0} \right)^T - \lambda_0^T \right] \delta_0 + \sum_{k=0}^{N-1} \left[\left(\frac{\partial H_k}{\partial \xi_k} \right)^T - \lambda_k^T \right] \delta_k + \sum_{k=0}^{N-1} \left(\frac{\partial H_k}{\partial y_k} \right)^T \eta_k = 0.$$

Since the indicated variations are mutually independent, following individual conditions must be satisfied

$$(90) \quad \lambda_k = \frac{\partial H_k}{\partial \xi_k},$$

$$(91) \quad \frac{\partial H_k}{\partial y_k} = 0$$

and the transversality conditions

$$(92) \quad \left[\left(\frac{\partial \Theta_N}{\partial \xi_N} \right)^T - \lambda_N^T \right] \delta_N = 0,$$

$$(93) \quad \left[\left(\frac{\partial \Theta_0}{\partial \xi_0} \right)^T - \lambda_0^T \right] \delta_0 = 0.$$

If the value of any variable is specified, the corresponding variation vanishes and the respective condition (90) through (93) does not apply. Particularly for given ξ_0 and ξ_N the corresponding boundary condition on λ_k is satisfied by δ_0 and δ_N respectively being both equal to zero.

The second condition (86) is satisfied for all cost functions and systems of interest.

For the linear regulator problem with specified ξ_0 and ξ_N eqs (90) and (91) yield

$$(94) \quad \frac{\partial H_k}{\partial \xi_k} = Q_k \xi_k + F_k^T \lambda_{k+1} = \lambda_k,$$

$$(95) \quad \frac{\partial H_k}{\partial y_k} = R_k y_k + G_k^T \lambda_{k+1} = 0.$$

These are Euler equations for the variational problem under consideration. The estimated solution of these equations is 99

$$(96) \quad \lambda_k = P_k \xi_k.$$

Combining now eqs. (95), (96) and (5), we have

$$(97) \quad R_k y_k + G_k^T P_{k+1} (F_k \xi_k + G_k y_k) = 0$$

and the control vector is

$$(98) \quad y_k = - (G^T P_{k+1} G_k + R_k)^{-1} G_k^T P_{k+1} F_k \xi_k.$$

This is exactly the same result as indicated by eq. (64) and consequently, eqs. (70) and (71) are valid for the maximum principle, too.

The Riccati equation can be derived by means of eq. (94) if eqs. (96), (5) and (98) are used simultaneously.

$$(99) \quad P_k \xi_k = Q_k \xi_k + F_k^T P_{k+1} [F_k \xi_k - G_k (G_k^T P_{k+1} G_k + R_k)^{-1} G_k^T P_{k+1} F_k \xi_k].$$

Eq. (99) must hold for any vector ξ_k . Applying this condition on eq. (99), after some rearrangements, the same form of the Riccati equation is obtained as introduced by eq. (72).

If for example ξ_N would not be specified, then an additional condition must be applied, i.e.

$$(100) \quad \lambda_N = K \xi_N$$

with K being nonnegative definite matrix in order for the second variation to be positive. Notice, that for the final stage of the process

$$(101) \quad P_N = K$$

must hold.

5.2. Stationary controlled systems for N approaching to infinity

For time invariant controlled systems and for constant weighting matrices of the cost function, relation (99) yields the Riccati equation of the form

$$(102) \quad P = F^T P F - F^T P G (G^T P G + R)^{-1} G^T P F + Q$$

corresponding to eq. (72) with all indexes omitted. The control vector is according to (98)

$$(103) \quad y_k = - (G^T P G + R)^{-1} G^T P F \xi_k.$$

It is obvious from eq. (88), that for fixed ξ_0 and ξ_N i.e. for $\delta_0 = \delta_N = 0$ and for N approaching to infinity, it holds that

$$(104) \quad \sum_{k=0}^{N-1} \lambda_{k+1}^T \delta_{k+1} = \sum_{k=0}^{N-1} \lambda_k^T \delta_k = \sum_{k=0}^N \lambda_k^T \delta_k.$$

Hence for this particular problem it is possible to change the index $k+1$ into k in eq. (85) and to modify next relations respectively. Notice, that the cost function is

$$(105) \quad J_\infty = \sum_{k=0}^{\infty} (H_k - \lambda_k^T \xi_{k+1})$$

and the Euler equations take the form

$$(106) \quad \frac{\partial H_k}{\partial \xi_k} = Q \xi_k + F^T \lambda_k = \lambda_{k-1},$$

$$(107) \quad \frac{\partial H_k}{\partial y_k} = R y_k + G^T \lambda_k = 0.$$

Using the same procedure for the derivation of the control vector as in the preceding section, we have

$$(108) \quad y_k = -R^{-1} G^T P \xi_k.$$

Evidently the weighting matrix R must be positive definite in order that the inverse R^{-1} exists. However this mathematical result is in discrepancy with the physical reality then, as it was mentioned in paragraph 4.1, for discrete linear systems it is not necessary to constrain the controlling variables. If such a constraint is expressed by matrix R it may concern only some controlling variables of the vector y . Nevertheless the relation (108) is correct for matrix R being positive definite i.e., from the physical point of view, for all controlling variables of the vector y constrained.

For the Riccati equation λ_k is to be calculated from eq. (106). With index k changed into $k+1$, we have

$$(109) \quad \lambda_{k+1} = (F^T)^{-1} \lambda_k - (F^T)^{-1} Q \xi_{k+1}.$$

Combining now (109), (96), (5) and (108), the final result is

$$(110) \quad PF - [(F^T)^{-1} + (F^T)^{-1} QGR^{-1}G^T]P - PGR^{-1}G^TP + (F^T)^{-1}QF = 0.$$

At the first sight it is not clear whether eq. (110) has the same solution as eq. (102). Moreover, the control vector (108) has quite another form than control vector (103). The question is whether the both mentioned solutions are identical.

To examine this problem, we can write the so-called matrix Euler equation

$$(111) \quad p_{k+1} = Ep_k$$

with vector

$$(112) \quad p_k = \begin{bmatrix} \xi_k \\ \lambda_k \end{bmatrix}.$$

If for both cases the eigenvalues of the Euler matrix E of the homogeneous equation (111) will be the same then both optimal solutions are identical.

Proof. For the set of eqs. (94) and (95) it holds

$$(113) \quad \xi_{k+1} = F\xi_k - GR^{-1}G^T\lambda_{k+1},$$

$$(114) \quad \lambda_{k+1} = -(F^T)^{-1} Q\xi_k + (F^T)^{-1} \lambda_k.$$

Substituting (114) for λ_{k+1} in (113), we obtain

$$(115) \quad \xi_{k+1} = [F + GR^{-1}G^T(F^T)^{-1}Q] \xi_k - GR^{-1}G^T(F^T)^{-1} \lambda_k.$$

Eqs. (115) and (114) define a system the Euler matrix of which has the form

$$(116) \quad E = \begin{bmatrix} F + GR^{-1}G^T(F^T)^{-1}Q & -GR^{-1}G^T(F^T)^{-1} \\ -(F^T)^{-1}Q & (F^T)^{-1} \end{bmatrix}.$$

To examine the eigenvalues of the matrix (116) we write the characteristic matrix $E - \lambda I$ with I being the identity matrix. The elements of the characteristic matrix are in general polynomials. For such a type of matrices the eigenvalues remain unchanged if the respective matrix is rearranged by the application of admissible changes which all may be expressed by nonsingular transform matrices. Hence, for the particular case, it holds

$$(117) \quad |E - \lambda I| = \left| \begin{bmatrix} I & GR^{-1}G^T \\ 0 & F \end{bmatrix} \begin{bmatrix} F + GR^{-1}G^T(F^T)^{-1}Q - \lambda I & -GR^{-1}G^T(F^T)^{-1} \\ -(F^T)^{-1}Q & (F^T)^{-1} - \lambda I \end{bmatrix} \right|.$$

$$= \left| \begin{bmatrix} I & 0 \\ 0 & -\frac{I}{\lambda} \end{bmatrix} \right|,$$

$$(118) \quad |E - \lambda I| = \begin{vmatrix} F - \lambda I & GR^{-1}G^T \\ -Q & F^T - \frac{I}{\lambda} \end{vmatrix}.$$

On the other hand, for the set of Euler eqs. (106) and (107) by the same procedure, we obtain

$$(119) \quad \xi_{k+1} = F\xi_k - GR^{-1}G^T\lambda_k,$$

$$(120) \quad \lambda_{k+1} = -(F^T)^{-1} Q\xi_{k+1} + (F^T)^{-1} \lambda_k.$$

102 Substituting from (119) ξ_{k+1} into (120) we have

$$(121) \quad \lambda_{k+1} = -(F^T)^{-1} Q F \xi_k + [(F^T)^{-1} Q G R^{-1} G^T + (F^T)^{-1}] \lambda_k.$$

Eqs. (119) and (121) define the Euler matrix

$$(122) \quad E = \begin{bmatrix} F & -GR^{-1}G^T \\ -(F^T)^{-1}QF & (F^T)^{-1}QGR^{-1}G^T + (F^T)^{-1} \end{bmatrix}.$$

For the characteristic matrix it holds

$$(123) \quad |E - \lambda I| = \begin{vmatrix} I & 0 \\ Q & F^T \end{vmatrix} \begin{vmatrix} F - \lambda I & -GR^{-1}G^T \\ -(F^T)^{-1}QF & (F^T)^{-1}QGR^{-1}G^T + (F^T)^{-1} - \lambda I \end{vmatrix}.$$

$$= \begin{vmatrix} I & 0 \\ 0 & -\frac{I}{\lambda} \end{vmatrix},$$

$$(124) \quad |E - \lambda I| = \begin{vmatrix} F - \lambda I & GR^{-1}G^T \\ -Q & F^T - \frac{I}{\lambda} \end{vmatrix}.$$

From the results (118) and (124) it is clear that for both cases under consideration the characteristic polynomials (or the characteristic matrices) of the Euler matrix have the same form and consequently they have the same eigenvalues and both solutions represent identically the same optimal system.

Remark. Evidently there is a question whether similar relations are valid for non-stationary systems. We shall leave to consider this question to diligent readers.

The equality of Euler matrices (118) and (124) does not prove that the solution of both respective Riccati equations (102) and (110) is identically the same. On the other hand we can assume that the both control vectors (103) and (108) respectively must yield the same values of the controlling variable y_k , $k = 1, 2, \dots$ if all other conditions are identical. A problem which is of considerable interest to us is to examine whether the equality of values of the controlling variable for both cases under consideration is reached by different matrices P satisfying the Riccati equations (102) and (110) respectively or whether the solutions of the mentioned Riccati equations are identically the same. This may be the subject of our next investigation.

Starting with the control vector of the form (103), we shall rearrange first the following relation

$$(125) \quad \begin{aligned} -M &= (G^T P G + R)^{-1} G^T P F = \\ &= R^{-1} (G^T P G + R - G^T P G) (G^T P G + R)^{-1} G^T P F \end{aligned}$$

where $R^{-1}(G^T P G + R - G^T P G)$ is identity matrix. Hence

$$(126) \quad \begin{aligned} -M &= R^{-1}[I - G^T P G(G^T P G + R)^{-1}] G^T P F = \\ &= R^{-1} G^T [P - P G(G^T P G + R)^{-1} G^T P] F. \end{aligned}$$

The term in the brackets corresponds according to the matrix lemma to $(P^{-1} + GR^{-1}G^T)^{-1}$. Finally we have

$$(127) \quad (G^T P G + R)^{-1} G^T P F = R^{-1} G^T (P^{-1} + GR^{-1}G^T)^{-1} F.$$

Denoting now the solution of the Riccati equations (102) by P_1 and that one of Riccati equation (110) by P_2 , it is evident that on the right hand side of eq. (126) we have the function $(-M)$ of the form (108) under the assumption that

$$(128) \quad P_2 = (P_1^{-1} + GR^{-1}G^T)^{-1} F.$$

Hence the solutions of both Riccati equations (102) and (110) respectively are different and eq. (128) represents the mutual relation.

Proof. To verify this result, it is possible to insert relation (128) into the equation (110) and the other form of the Riccati equation must be obtained.

To simplify this procedure it is useful to modify eq. (102) and (110) into

$$(129) \quad P_1 F - (F^T)^{-1} P_1 - P_1 G(G^T P_1 G + R)^{-1} G^T P_1 F + (F^T)^{-1} Q = 0$$

and

$$(130) \quad P_2 F - (F^T)^{-1} P_2 - (F^T)^{-1} Q G R^{-1} G^T P_2 - P_2 G R^{-1} G^T P_2 + (F^T)^{-1} Q F = 0$$

respectively.

The third term in (129) may be expressed as

$$(131) \quad P_1 G(G^T P_1 G + R)^{-1} G^T P_1 F = P_1 F - (F^T)^{-1} P_1 + (F^T)^{-1} Q$$

which is the last term of the right hand side of eq. (126). Combining now relations (126), (127) and (131), we obtain for the control vector

$$(132) \quad R^{-1} G^T [P_1 F - P_1 F + (F^T)^{-1} P_1 - (F^T)^{-1} Q] = R^{-1} G^T (F^T)^{-1} (P_1 - Q).$$

Comparing the final result of (132) with the right hand side of (127) it is evident that

$$(133) \quad P_2 = (F^T)^{-1} (P_1 - Q).$$

This relation is more suitable for substitution into Riccati equation (130) than the

104 previous one, eq. (128). Using the latter form we obtain

$$(134) \quad \begin{aligned} & (F^T)^{-1} (P_1 - Q) F - (F^T)^{-2} (P_1 - Q) - \\ & - (F^T)^{-1} Q G R^{-1} G^T (F^T)^{-1} (P_1 - Q) - \\ & - (F^T)^{-1} (P_1 - Q) G R^{-1} G^T (F^T)^{-1} (P_1 - Q) + (F^T)^{-1} Q F = 0, \\ & P_1 F - (F^T)^{-1} P_1 + (F^T)^{-1} Q - P_1 G R^{-1} G^T (F^T)^{-1} (P_1 - Q) = 0. \end{aligned}$$

Since

$$(135) \quad R^{-1} G^T (F^T)^{-1} (P_1 - Q) = (G^T P_1 G + R)^{-1} G^T P_1 F$$

eq. (134) can be easily modified into the Riccati equation of the form (129), which was to be proved. At this opportunity it is worth to mention that the solution of the Riccati equation of the form (110) is somewhat easier than that one of the form (102). However the condition relating to the matrix R to be positive definite represents a very strong limitation.

6. DYNAMIC PROGRAMMING

6.1 General procedure

Assuming again the same problem as in section 4 and 5, the optimum digital control problem may be considered as N -stage decision process. Using for the determination of the optimum the dynamic programming of Bellman, we obtain necessarily identically the same results as in preceding paragraphs. The procedure of dynamic programming is very well known and described several times in the technical literature. Hence we shall pay attention only to the main steps of the procedure enabling us to derive the desired results.

The minimum value of the performance index can be denoted by

$$(136) \quad \begin{aligned} f_{0,N} &= f_N - f_0 = \\ &= \xi_N^T P_N \xi_N + \min_{y_i} \sum_{i=0}^{N-1} [\xi_i^T Q_i \xi_i + y_i^T R_i y_i]. \end{aligned}$$

This form of the cost function corresponds to that one given by eq. (18). A more general form can be expressed in the following way

$$(137) \quad f_{k,N} = \xi_N^T P_N \xi_N + \min_{y_j} \sum_{j=k}^{N-1} [\xi_j^T Q_j \xi_j + y_j^T R_j y_j] = \xi_N^T P_N \xi_N + \min_{y_j} J_{k,N}$$

with $j = k - 1, k, \dots, N - 1$ and for $k = 0, 1, 2, \dots, N - 1$. For $k = 0$ relation (137) reduces to (136) and the first vector ξ_0 is given by initial conditions.

Assuming that the value of the cost function corresponding to the first $k - 1$ stages is optimum, then the increase of this value in the remaining $N - k$ stages is equal to the increment corresponding to the stage k plus optimum increases in the next $N - (k + 1)$ stages. Hence the optimum value of the cost function in the $N - k$ stages is

$$(138) \quad f_{k,N} = \xi_N^T P_N \xi_N + \min_{y_j} [\xi_k^T Q_k \xi_k + y_k^T R_k y_k + f_{k+1,N}]$$

with $j = k - 1, k, \dots, N - 1$.

Since the cost function and the functional f are quadratic in ξ , it can be expected that

$$(139) \quad f_{k,N} = \xi_k^T P_k \xi_k$$

for $k = 0, 1, 2, \dots, N$. With relation (139), eq. (138) becomes

$$(140) \quad f_{k,N} = \xi_N^T P_N \xi_N + \min_{y_j} [\xi_k^T Q_k \xi_k + y_k^T R_k y_k + \xi_{k+1}^T P_{k+1} \xi_{k+1}] = \\ = \xi_N^T P_N \xi_N + \min_{y_j} J_{k,N}, \\ j = k - 1, k, \dots, N - 1.$$

Substituting for ξ_{k+1} according to the state equation, we have

$$(141) \quad J_{k,N} = [\xi_k^T Q_k \xi_k + y_k^T R_k y_k + (\xi_k^T F_k^T + y_k^T G_k^T) P_{k+1} (G_k y_k + F_k \xi_k)].$$

Differentiating the last relation with respect to y_k yields

$$(142) \quad \frac{\partial J_{k,N}}{\partial y_k} = 2(R_k y_k + G_k^T P_{k+1} G_k) y_k + 2G_k^T P_{k+1} F_k \xi_k.$$

At the minimum the derivative is zero and thus the control vector is

$$(143) \quad y_k = M_k \xi_k = -(G_k^T P_{k+1} G_k + R_k)^{-1} G_k^T P_{k+1} F_k \xi_k.$$

The respective Riccati equation can be derived by means of relation (141) when (143) is inserted for y_k and $J_{k,N}$ is expressed according to (139) as $\xi_k^T P_k \xi_k$. We obtain

$$(144) \quad \xi_k^T P_k \xi_k = \xi_k^T Q_k \xi_k + \xi_k^T F_k^T P_{k+1} G_k (G_k^T P_{k+1} G_k + R_k)^{-1} \cdot \\ \cdot R_k (G_k^T P_{k+1} G_k + R_k)^{-1} G_k^T P_{k+1} F_k \xi_k + \\ + [\xi_k^T F_k^T - \xi_k^T F_k^T P_{k+1} G_k (G_k^T P_{k+1} G_k + R_k)^{-1} G_k^T] \cdot \\ \cdot P_{k+1} [-G_k (G_k^T P_{k+1} G_k + R_k)^{-1} G_k^T P_{k+1} F_k \xi_k + F_k \xi_k].$$

106 Eq. (144) must be valid for arbitrary ξ_k and consequently it holds that

$$(145) \quad \begin{aligned} P_k &= Q_k + F_k^T P_{k+1} G_k (G_k^T P_{k+1} G_k + R_k)^{-1} \cdot \\ &\cdot [R_k (G_k^T P_{k+1} G_k + R_k)^{-1} - 2 + G_k^T P_{k+1} G_k \cdot \\ &\cdot (G_k^T P_{k+1} G_k + R_k)^{-1}] G_k^T P_{k+1} F_k + F_k^T P_{k+1} F_k. \end{aligned}$$

Finally we have

$$(146) \quad P_k = Q_k - F_k^T P_{k+1} G_k (G_k^T P_{k+1} G_k + R_k)^{-1} G_k^T P_{k+1} F_k + F_k^T P_{k+1} F_k.$$

This is again the same equation as previously derived.

6.2 Special case

We shall consider only one special problem i.e. the infinite-stage process when $N \rightarrow \infty$. In this case eq. (138) reduces to

$$(147) \quad f_{k,\infty} = \min_{y_j} [\xi_k^T Q_k \xi_k + y_k^T R_k y_k + f_{k+1,\infty}],$$

$$j = k - 1, k, \dots, \infty$$

and relation (139) may be written as

$$(148) \quad f_{k,\infty} = \xi_k^T P \xi_k$$

where P is a constant matrix. By the same procedure as described above, it is possible to obtain the modified result with $P_{k+1} = P_k = P$. Since the Riccati equation (146) must be satisfied for any k and the matrix P is a constant matrix, the weighting matrices Q_k and R_k cannot be selected arbitrarily but with respect to the non-stationarity of the controlled system in such a way to satisfy the respective form of the Riccati equation in order to achieve the optimum control. For example if P once known, it is possible to select R_k for any k and calculate Q_k satisfying the Riccati equation.

All results corresponding to stationary controlled systems and cost functions not time weighted follow directly from the more common results given in section 6.1.

7. SOLUTION OF THE STATIONARY PROBLEM

In previous sections we derived relations necessary for the solution of the optimum process but the solution itself was not described. In most modifications the question is to find the solution of the Riccati equation and knowing once the matrix P it is possible to calculate the control vector. Understand, in the continuous version, usually it is assumed to calculate the matrix P directly by a suitable numerical iterative procedure i.e. without calculating the eigenvalues of the characteristic

matrices $(K - \lambda I)$ or $(E - \lambda I)$ where K is a transition matrix of the closed control loop and E an Euler matrix respectively. These procedures possess without any doubt a significant practical importance. However, since these procedures represent purely numerical algorithms, they deviate from the aim of this article and therefore it is assumed to describe the appropriate methods in a special paper.

On the other hand for the discrete version solution of the "Riccati equation" is much simpler even in the nonstationary case.

It is possible to state that at the end point of the control process i.e. in the stage N , the increment of the cost function is just $\xi_N^T Q_N \xi_N$ then the component corresponding to the controlling variable does not apply. This increment is e.g. according to (139) $\xi_N^T P_N \xi_N$. Hence for $k = N$ it holds that $P_N = Q_N$. Knowing P_N , we can calculate the feedback transition matrix M_{N-1} using (143), then P_{N-1} from the equation (146), next M_{N-2} , then P_{N-2} etc.

Now we shall draw our attention to the more or less classical method using the knowledge of the eigenvalues of Euler matrix. This method is described here for stationary systems and not time varying weighting matrices of the cost function and it applies an obvious cancelation of unstable components of the solution and has a close relation to the Riccati equation.

The general procedure described in this section may be applied for the set of equations (94) and (95) or (106) and (107) respectively. For the sake of brevity only the first set of Euler equations will be used for the demonstration of the method.

Theorem 5. *The eigenvalues of the Euler matrix E (eqs. (116) and (122)) are symmetrically displaced with respect to the unit circle centred in the origin of the complex plane in that sense that from the total of $2n$ eigenvalues n of the value $|\lambda_i| < 1$, $i = 1, 2, \dots, n$ are stable and n of the value $|\lambda_i^{-1}| > 1$ are unstable provided that there are no eigenvalues $\lambda_i = 1$.*

Proof. It holds for the polynomial characteristic matrix that the determinant is unchanged if the matrix is transposed. According to (118) we may write

$$(149) \quad \begin{vmatrix} F - \lambda I & GR^{-1}G^T \\ -Q & F^T - \lambda^{-1}I \end{vmatrix} = \begin{vmatrix} F^T - \lambda I & -Q \\ GR^{-1}G^T & F - \lambda^{-1}I \end{vmatrix}.$$

Interchanging an odd number of rows or columns, the sign is changed. Hence

$$(150) \quad \begin{vmatrix} F^T - \lambda I & -Q \\ GR^{-1}G^T & F - \lambda^{-1}I \end{vmatrix} = - \begin{vmatrix} GR^{-1}G^T & F - \lambda^{-1}I \\ F^T - \lambda I & -Q \end{vmatrix} = \begin{vmatrix} F - \lambda^{-1}I & GR^{-1}G^T \\ -Q & F^T - \lambda I \end{vmatrix}.$$

Comparing the first term in (149) and the last term in (150) it is evident that the characteristic matrix is invariant with respect to change of λ into λ^{-1} and consequently theorem 5 is proved.

Assuming now the system described by the homogeneous equation (111), it is

108 possible to write the general form of the solution

$$(151) \quad p_k = E^k p_0 = T J_E^k T^{-1} p_0$$

where J_E is a Jordan matrix of E and T is a matrix the columns of which represent a complete set of eigenvectors of E . Applying theorem 5, it is advantageous to group the eigenvalues and the corresponding eigenvectors of E in such a way that the n stable eigenvalues and the corresponding eigenvectors come first. Accordingly it is possible to introduce following denotations

$$(152) \quad J_E = \begin{bmatrix} J_{11} & 0 \\ 0 & J_{22} \end{bmatrix},$$

$$(153) \quad T = \begin{bmatrix} T_{11} & T_{12} \\ T_{21} & T_{22} \end{bmatrix},$$

$$(154) \quad T^{-1} = T^* = \begin{bmatrix} T_{11}^* & T_{12}^* \\ T_{21}^* & T_{22}^* \end{bmatrix}$$

where all indexed matrices are of the dimension $(n; n)$ and matrices E , J_E , T of the dimension $(2n; 2n)$. Hence eq. (151) can be rewritten into

$$(155) \quad \begin{bmatrix} \xi_k \\ \lambda_k \end{bmatrix} = \begin{bmatrix} T_{11} & T_{12} \\ T_{21} & T_{22} \end{bmatrix} \begin{bmatrix} J_{11}^k & 0 \\ 0 & J_{22}^k \end{bmatrix} \begin{bmatrix} T_{11}^* & T_{12}^* \\ T_{21}^* & T_{22}^* \end{bmatrix} \begin{bmatrix} \xi_0 \\ \lambda_0 \end{bmatrix}$$

and defines two sets of equations

$$(156) \quad \xi_k = (T_{11} J_{11}^k T_{11}^* + T_{12} J_{22}^k T_{21}^*) \xi_0 + (T_{11} J_{11}^k T_{12}^* + T_{12} J_{22}^k T_{22}^*) \lambda_0,$$

$$(157) \quad \lambda_k = (T_{21} J_{11}^k T_{11}^* + T_{22} J_{22}^k T_{21}^*) \xi_0 + (T_{21} J_{11}^k T_{12}^* + T_{22} J_{22}^k T_{22}^*) \lambda_0.$$

Since stable solution is required, all terms in eq. (156) multiplied by the unstable field J_{22}^* must vanish. Therefore

$$(158) \quad T_{21}^* \xi_0 + T_{22}^* \lambda_0 = 0$$

which yields the relation for λ_0

$$(159) \quad \lambda_0 = -(T_{22}^*)^{-1} T_{21}^* \xi_0$$

provided that initial conditions ξ_0 are known. On the other hand, with respect to the relation (96), it holds that

$$(160) \quad P = -(T_{22}^*)^{-1} T_{21}^*.$$

Eq. (156) yields the vector of state variables

$$(161) \quad \xi_k = T_{11} J_{11}^k (T_{11}^* \xi_0 + T_{12}^* \lambda_0).$$

Substituting (159) for λ_0 into (161) we have

$$(162) \quad \xi_k = T_{11} J_{11}^k [T_{11}^* - T_{12}^* (T_{22}^*)^{-1} T_{21}^*] \xi_0.$$

By means of (160) and (162) it is possible to calculate the control vector y according to (103).

Besides by eq. (160) matrix P may be expressed by

$$(163) \quad P = T_{21} (T_{11})^{-1}$$

too if for (96) λ_k is calculated from (157) and ξ_k is substituted in accordance with (162). Notice that condition (158) is valid again. It is easy to prove that (160) and (163) yield identically the same values.

Remark. In section 5.2 was proved that the eigenvalues of Euler matrices (116) and (122) are identically the same. On the other hand, relations (160) or (163) must yield different results for each of the mentioned Euler matrices in order to satisfy eq. (128). Consequently the eigenvectors of matrices (116) and (122) must be different.

Theorem 6. *If R is a positive definite matrix then eigenvalues of the Euler matrix (122) are eigenvalues of the matrix of the closed control loop too.*

Proof. Substituting the control vector (108) into the state equation (5) it is possible to derive the transition matrix of the closed control loop

$$(164) \quad K = F - GR^{-1}G^T P$$

which substituted into the Riccati equation (130) yields

$$(165) \quad PK = [(F^T)^{-1} + (F^T)^{-1} QGR^{-1}G^T] P - (F^T)^{-1} QF.$$

Let the matrix V transform the matrix K into the Jordan form, so that

$$(166) \quad K = VJ_K V^{-1}$$

and let according to (163)

$$(167) \quad P = UV^{-1}.$$

Then, inserting (166) and (167) into (164) and (165), we obtain

$$(168) \quad VJ_K = FV - GR^{-1}G^T U,$$

$$(169) \quad UJ_K = [(F^T)^{-1} + (F^T)^{-1} QGR^{-1}G^T] U - (F^T)^{-1} QFV.$$

or

$$(170) \quad \begin{bmatrix} V \\ U \end{bmatrix} J_K = E \begin{bmatrix} V \\ U \end{bmatrix}; \quad \begin{bmatrix} V \\ U \end{bmatrix} = N.$$

110 Let a_1, a_2, \dots, a_n be the columns of the matrix N of the dimension $2n \times n$ and λ_i , $i = 1, 2, \dots, n$ be the eigenvalues of the matrix K . For distinct eigenvalues λ_i it holds according to (170)

$$(171) \quad a_i \lambda_i = E a_i, \quad i = 1, 2, \dots, n.$$

Hence λ_i are eigenvalues of E too and a_i are its corresponding eigenvectors which was to be proved. Since λ_i are assumed to be stable for the closed control loop, matrices U and V in (167) must correspond to T_{21} and T_{11} in (163) respectively.

Similarly for multiple λ_i

$$(172) \quad \begin{aligned} a_i \lambda_i &= E a_i, \\ a_i + a_{i+1} \lambda_i &= E a_{i+1}, \\ &\vdots \\ a_{i+p_i-2} + a_{i+p_i-1} \lambda_i &= E a_{i+p_i-1}, \end{aligned}$$

where p_i signifies the dimension of the Jordan field corresponding to eigenvalue λ_i . Since V as a transforming matrix is nonsingular, $a_i \neq 0$, λ_i is an eigenvalue of the matrix E too and $a_i, a_{i+1}, \dots, a_{i+p_i-1}$ is the respective chain of the length p_i of generalized eigenvectors. The solution may be expressed again by (163) or (167).

Note. For R positive semidefinite it is possible to find the solution according to the general procedure described in this section and relating to the Euler matrix (116). On the other hand it is possible to modify the procedure described in this section for R positive definite.

Assume Q and R to be positive semidefinite and symmetric matrices. The rank of the matrix Q is p . Introduce

$$(173) \quad Q = C^*{}^T C^*$$

where the matrix C^* has dimension $(p; n)$ and the rank p .

Let be

$$(174) \quad S = C^*{}^T \tilde{S}.$$

Denote

$$\tilde{x}_k = C^* \zeta_k$$

and assume this variable \tilde{x}_k in accordance with eq. (6) as an output variable of some controlled plant with shifted output, $m > 0$. Notice that this fictitious controlled plant can differ from the given plant, so far as the matrix C^* defined by (173) differs from the output transition matrix C of the given plant or if the given plant does not have shifted output, i.e. $m = 0$.

With (173), (174) and (175) it is possible to write the cost function in a new form

111

$$(176) \quad J = \sum_{k=0}^{\infty} (\tilde{x}_k^T \tilde{x}_k + 2\tilde{x}_k^T \tilde{S} y_k + y_k^T R y_k).$$

Since the fictitious plant has shifted output by $m > 0$ periods of sampling, the output values $\tilde{x}_0, \tilde{x}_1, \dots, \tilde{x}_{m-1}$ are not influenced by the controlling variable y_k , $m-1 \geq k \geq 0$ and consequently it is not necessary to have these values in the cost function. Hence

$$(177) \quad J^* = \sum_{k=0}^{\infty} (\tilde{x}_{k+m}^T \tilde{x}_{k+m} + 2\tilde{x}_{k+m}^T \tilde{S} y_k + y_k^T R y_k)$$

and the optimal control according to (177) will be identical with that one obtained when using the cost function (176).

Now there arises a problem to calculate \tilde{x}_{k+m} . Since \tilde{x}_k is not influenced by $y_k, y_{k-1}, \dots, y_{k-m+1}$, \tilde{x}_k can be expressed in accordance with (27) as

$$(178) \quad \tilde{x}_k = C^* F^k \xi_0 + C^* F^{k-1} G y_0 + \dots + C^* F^{m-1} G y_{k-m}.$$

Now, y_{k-m} calculated from (178) is

$$(179) \quad y_{k-m} = (C^* F^{m-1} G)^{-1} \tilde{x}_k - (C^* F^{m-1} G)^{-1} [C^* F^k \xi_0 + C^* F^{k-1} G y_0 + \dots + C^* F^m G y_{k-m-1}].$$

Notice, that the expression in the brackets equals in accordance with (26) to $C^* F^m \xi_{k-m}$. Hence

$$(180) \quad y_{k-m} = \Gamma_m^{-1} \tilde{x}_k - \Gamma_m^{-1} C^* F^m \xi_{k-m}$$

where, with the abbreviated notation from theorem 2,

$$(181) \quad \Gamma_m = C^* F^{m-1} G.$$

Eq. (180) is valid for any k , consequently also in the form

$$(182) \quad y_k = \Gamma_m^{-1} \tilde{x}_{k+m} - \Gamma_m^{-1} C^* F^m \xi_k$$

which yields the desired state vector

$$(183) \quad \tilde{x}_{k+m} = C^* F^m \xi_k + \Gamma_m y_k.$$

Substituting (183) into (177) we obtain

$$(184) \quad J^* = \sum_{k=0}^{\infty} (\xi_k^T Q^* \xi_k + 2\xi_k^T S^* y_k + y_k^T R^* y_k)$$

$$\begin{aligned}
 (185) \quad Q^* &= (F^T)^m C^{*T} C F^m, \\
 S^* &= (F^T)^m C^{*T} \Gamma_m + S, \\
 R^* &= \Gamma_m^T \Gamma_m + R.
 \end{aligned}$$

By means of relations (21) we can eliminate in (184) the term with matrix S^* and obtain an usual form of the cost function. Solution of the given problem can be performed by an appropriate procedure described in this paper. From the last relation in (185) is evident that matrix R can equal to zero as well.

In conclusion it is possible to state that the problem of optimal control with cost function (18), where matrices Q and R are positive semidefinite and symmetric and where the rank of the matrix Q is p , can be transformed into an equivalent problem with cost function (184). There is R^* a symmetric matrix if the matrix Γ_m is positive definite and if its rank is $\gamma = r \leq p$. The procedure described in this note is valid for N being a finite integer or for $N \rightarrow \infty$ as well.

Acknowledgement. The author wishes to thank Mrs. Ing. A. Halousková PhD for her reviewing the manuscript and for most helpful and constructive comments improving the presentation of the material.

(Received September 14, 1971.)

REFERENCES

- Fan L. T. and Wang C. S.: The Discrete Maximum Principle. John Wiley & Sons, New York 1964.
 Gantmacher F. R.: Theory of Matrices. Vols. I and II, Chelsea Publishing Co., New York 1959.
 Kalman R. E.: Mathematical Description of Linear Dynamical Systems. J.S.I.A.M. Control, Ser. A, 1 (1963), 2, 159–192.
 Kučera V.: Optimal Control of Linear Discrete Systems According to Quadratic Cost Function (in Czech. Language). Research Report, ÚTIA ČSAV, Prague 1971.
 Ogata K.: State space Analysis of Control Systems. Prentice-Hall, Inc., Englewood Cliffs, N.Y. 1967.
 Pearson J. B. and Sridhar R.: A Discrete Optimal Control Problem. IEEE Transactions on Automatic Control AC-11, (April, 1966), 2.
 Pontryagin L. S. et al.: The Mathematical Theory of Optimal Processes. Interscience Publishers, New York 1962.
 Sage A. P.: Optimum Systems Control. Prentice-Hall, Inc., Englewood Cliffs, N.Y. 1968.
 Strejc V.: State Space Equations for Control Theory (in Czech Language). Supplement of the journal Kybernetika 5 (1969).
 Tou J. T.: Optimum Design of Digital Control Systems. Academic Press, New York, London 1963.
 Zadeh L. A. and Desoer C. A.: Linear System Theory, The State Space Approach. McGraw-Hill, 1963.

Syntéza diskrétních lineárních systémů ve stavovém prostoru

VLADIMÍR STREJC

Cílem tohoto článku je upozornit na některé důležité souvislosti mezi různými postupy syntézy ve stavovém prostoru diskrétních lineárních systémů optimalizuje-li se kvadratické kritérium jakosti. Obecná forma syntézy se vztahuje k jednorozměrným a mnohorozměrným úlohám řízení řešeným pomocí druhé Ljapunovovy věty, principu maxima Pontrjagina a pomocí Bellmanova dynamického programování. Je snahou předložit zhuštěný souhrn, který zdůrazňuje základní principy a integruje osvědčené postupy na ucelený obraz.

Prof. Ing. Vladimír Strejc, DrSc., Ústav teorie informace a automatizace ČSAV (Institute of Information Theory and Automation — Czechoslovak Academy of Sciences), Vítězslavská 49, Praha 2.