

Vladimir Semenovich Pugachev
Optimal learning systems

Kybernetika, Vol. 7 (1971), No. 5, (347)--376

Persistent URL: <http://dml.cz/dmlcz/125145>

Terms of use:

© Institute of Information Theory and Automation AS CR, 1971

Institute of Mathematics of the Academy of Sciences of the Czech Republic provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This paper has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library*
<http://project.dml.cz>

Optimal Learning Systems

V. S. PUGACHEV

Optimal Bayes learning systems and the main properties of their algorithms are studied from the general standpoint of Learning System Theory. The general concepts of learning system, teacher, forms of learning, types of teachers, etc. are discussed. The general case of a real teacher which performs control processes with random errors is considered. The conventional case of an ideal teacher which has been the only case previously considered, is treated as a special case. A measure of closeness of the performance of a learning system to the optimal system with complete information is introduced. Studying of optimal learning systems enables the designer of learning systems to determine extreme possibilities of learning systems and to estimate the performance of various specific algorithms of learning.

1. INTRODUCTION

A *learning system* is called such a system which improves itself using the information contained in signals received by it.

The main feature of a learning system distinguishing it from conventional automatic systems is the dependence of its output not on the actual input, only, but also on some previously obtained signals. Let Z be the input to the system, \hat{W} its output, \mathcal{E}' the set of previously obtained signals (teaching information) on which the output \hat{W} depends. Then the input-output relation of the learning system has the general form

$$(1.1) \quad \hat{W} = A(\mathcal{E}') Z,$$

where $A(\mathcal{E}')$ is some operator, deterministic or random, depending on teaching information \mathcal{E}' . For usual automatic systems the operator A is completely independent of previously obtained signals.

It should be emphasized that the notion of learning system concerns only systems improving themselves *while performing their functions*. All the information introduced into the system during its designing before it begins performing its functions can not be considered as teaching information. This information can only be treated as characterizing the initial organization of the system, its *capacities* to learn.

The period of time when the system receives the teaching signals E' on the basis of which its operator is formed, is called *the period of learning*. The period of learning may, in particular, include all the time of system working, i.e. the system may be permanently learning during all its life-time.

If the system does not receive additional external information apart its usual inputs during the period of learning, it is called *self-learning system*.

If the system receives some additional information as well as its usual inputs while learning, we say that it is *taught by a teacher*. A source of the additional information, being a man-operator or another automatic system, is called *teacher*.

The teacher may teach a system in two different ways. The first way is to show the system how to perform its functions. This way we shall call *teaching by show*. The second way is to observe the actions of the system and introduce into it teacher's estimates of its performance. We shall call this mode of teaching *teaching by estimating system actions*.

These three main ways of teaching may be combined sequentially or in parallel. For example, the system may be taught at first by show, then by estimating its actions, and then may continue improving its performance by self-learning. The system with several outputs may be taught simultaneously by show at some of its outputs, by estimating its actions at other outputs and may be self-learning with respect to the remaining outputs.

The teacher may know exactly the desired output W corresponding to a given input Z and show this desired output to the system. We shall call such a teacher *the ideal teacher*. Only the ideal teacher has been considered in all previous works in the field of learning systems. But the teacher does not in general know the exact desired output W and can only elaborate more or less suitable estimates of W . For example, the instructor teaching a man to pilot an aeroplane can not pilot this plane with absolute accuracy. He will inevitably make random errors resulting in random deviations of the plane from the required state of motion. In other words, he will show to the pupil only his own output, but can not show the actual desired output.

We shall call a teacher whose output \tilde{W} does not coincide with the desired output W *the real teacher*. The notion of the real teacher seems to be primarily introduced in the works of the author [1-3].

The algorithm of any real teacher is practically always stochastic, i.e. the output of the real teacher is random for any given input. So the algorithm of the real teacher can be defined only in the form of the conditional probability distribution of teacher output \tilde{W} given the input Z . We shall call this conditional distribution *decision function* of the teacher.

The necessity of teaching a system arises only if the probability distribution of the input and desired output or some parameters of this distribution are unknown. E.g., the distribution of feedback signals is unknown, if the dynamical characteristics of the controlled system are unknown. But the feedback signals represent some of the

components of the input to the control system. So the control system must learn to master performing the functions for which it is designed in this case.

We shall call *algorithm of learning* the mode of treating teaching signals to improve system performance. Any algorithm of learning which provides the increase of system performance is admissible in principle.

Up to now only heuristic algorithms of learning based on designer's skill and intuition were practically employed in learning system design.

The most natural approach to design algorithms of learning seems to be statistical estimation of unknown probability distribution of signals involved. This approach was initiated by Špaček and developed by his pupils [4–18] who studied some general properties of learning processes. In particular, the methods of stochastic approximations were extended and first used to derive the algorithms of learning [6, 8, 10–13]. The use of the methods of stochastic approximations to find algorithms of self-adaptive systems was also first proposed [15–17] (see also [19]). Following Špaček ideas the author proposed in [20] one possible statistical algorithm of self-learning for systems based on optimal Bayes decisions. Some special statistical problems related to self-learning were solved in [21–25]. The algorithms of learning for pattern recognition systems based on stochastic approximation methods were first derived for the case of the ideal teacher in [26] (see also [27]).

We shall derive here the general optimal algorithms of learning using Bayes statistical decision approach.* The system using such an optimal algorithm of learning will be called *Bayes optimal learning system*. This system possesses the best possible performance among all the learning systems receiving the same teaching information.

The study of Bayes optimal learning systems enables us to establish the extreme performance of learning systems, i.e. *potential learnability* of learning systems. Comparing the performance of any projected learning system with the performance of the Bayes optimal learning system one can estimate the algorithm of learning employed in the designed system.

We use here the terminology of automatic control theory, since the problems of teaching automatic systems are of special interest for us. Alternately, using the terminology accepted in studies of biological aspects of learning, Z represents the stimulus, W the desired reaction, \hat{W} the actual reaction of a learning system, \tilde{W} the consequence caused by system reaction \hat{W} to the stimulus Z .

2. BAYES OPTIMAL LEARNING SYSTEMS

We shall now give a precise formulation of the problem of finding a Bayes optimal learning system.

We suppose that all possible values z of the input Z represent elements of some

* This approach seems to be initiated by M. E. Shaikin [28] who solved one special problem of learning of pattern recognition systems.

set A , and all possible values w of the desired output W represent the elements of another set B . Let $\gamma(\Delta z; \theta | \lambda)$ be the family of probability measures on A , and $\kappa(\Delta w; \theta | z, \lambda)$ the family of conditional probability measures on B given the value z of Z , λ being a parameter with values in a set L , and θ a parameter whose values are real numbers*. Suppose that the input Z and desired output W represent random variables in A and B respectively distributed in accordance with the probability measure $\gamma(\Delta z; \theta | \lambda)$ and conditional probability measure $\kappa(\Delta w; \theta | z, \lambda)$ corresponding to some specific value of λ equal to the value $\lambda(\theta)$ of a random function $A(\theta)$ with values in L , the parameter θ taking different values in different cycles of the actions of the system (it may represent, in particular, the number of the cycle of system action).

If the true value of λ were known and introduced into the system at each cycle of its action, the system could make various statistical decisions, resulting in respective estimates of W . In particular, the system could evaluate the Bayes optimal estimate W_λ^* of W corresponding to any given loss function $l(W; \tilde{W} | \lambda)$ which may, in general, depend on the parameter λ . We shall call the system whose output at each cycle represents the Bayes optimal estimate W_λ^* of the desired output W corresponding to the true value of λ *Bayes optimal system with complete information about λ* .

If the true value of λ is unknown, the necessary estimates can be made, in principle, on the basis of some observations. Thus the system must *learn* to elaborate estimates of the desired output by making some observations and processing the information thus obtained, before it will begin performing the functions for which it is designed. We shall suppose that the learning period during which the system makes observations consists of N cycles of system action corresponding to values $\theta_1, \dots, \theta_N$ of the numerical parameter θ . The value of θ at the first cycle following the period of learning we shall leave without any indices.

If the system observes while learning only the inputs Z_1, \dots, Z_N , it is *self-learning*. If it receives, apart from the inputs Z_1, \dots, Z_N , some other signals, it is *learning with a teacher*. The teacher may show to the system estimates $\tilde{W}_1, \dots, \tilde{W}_N$ of the desired outputs W_1, \dots, W_N corresponding to the inputs Z_1, \dots, Z_N or estimates of some functionals of the pairs $(Z_1, W_1), \dots, (Z_N, W_N)$. This is the case of *teaching by show*. If the teacher observes the responses $\tilde{W}_1, \dots, \tilde{W}_N$ of the system corresponding to the inputs Z_1, \dots, Z_N and shows to the system estimates $\tilde{W}_1, \dots, \tilde{W}_N$ of its performance at each cycle, the system is *taught by estimating its actions*.

Thus the teacher output \tilde{W} may have, in general, another nature than the desired and actual system outputs W, \tilde{W} . Accordingly we shall suppose that all possible values \tilde{w} of the teacher output \tilde{W} represent elements of a set \tilde{B} which may, certainly,

* Speaking about measures defined on a set, we mean, certainly, measures defined on properly determined σ -algebra of subsets of this set. This will be assumed throughout the paper without further mentioning. We denote random variables by capital letters and their possible values (samples) by respective small letters; the sets of elements z, w, \dots we denote $\Delta z, \Delta w, \dots$, so Δz represents any subset of the set A , Δw any subset of the set B , etc.

coincide with B in the special case of teaching by show. The algorithm of the teacher is determined in general by the conditional probability measure $\delta_T(\Delta\tilde{w}; \theta | z, w, \hat{w}, \lambda)$ on \tilde{B} given the values z, w, \hat{w} of the input Z , desired output \tilde{W} , and actual system output \tilde{W} , λ being equal to the same value $\lambda(\theta)$ of the random function $A(\theta)$ as before. In the special case of the ideal teacher δ_T is condensed in single point w and does not depend on z, \hat{w}, λ . In the case of the real teacher teaching the system by show δ_T is independent of w, \hat{w} . In the case of the real teacher teaching the system by estimating its actions δ_T is independent of w .

We shall denote for brevity by Ξ the set of all signals received by the learning system including all the signals received during the learning period and the input Z received at the first cycle after learning.

In the case of self-learning Ξ represents the set of teaching inputs Z_1, \dots, Z_N followed by the input Z to which the system must respond in an optimal way. The probability measure of Ξ in the Cartesian product $X = A^{N+1}$ is determined in this case by

$$(2.1) \quad \sigma(\Delta\xi | \lambda, \tilde{\lambda}) = \gamma(\Delta z; \theta | \lambda) \prod_{i=1}^N \gamma(\Delta z_i; \theta_i | \lambda_i),$$

$\tilde{\lambda}$ being the set of values $\lambda_1, \dots, \lambda_N$ of the random function $A(\theta)$ corresponding to the values $\theta_1, \dots, \theta_N$ of θ . In deriving (2.1) the assumption was made that Z_1, \dots, Z_N, Z are conditionally independent, i.e. independent for any given set of values of $\lambda_1, \dots, \lambda_N, \lambda$. The measure $\sigma(\Delta\xi | \lambda, \tilde{\lambda})$ represents the conditional probability measure of Ξ given the values $\lambda_1, \dots, \lambda_N, \lambda$ of $A(\theta_1), \dots, A(\theta_N), A(\theta)$.

In the case of a system taught by show Ξ represents the set of previously obtained inputs Z_1, \dots, Z_N , corresponding teacher outputs $\tilde{W}_1, \dots, \tilde{W}_N$, and the input Z to which the system must respond optimally. The conditional probability measure of Ξ in the Cartesian product $X = A^{N+1} \times \tilde{B}^N$ is given in this case by

$$(2.2) \quad \sigma(\Delta\xi | \lambda, \tilde{\lambda}) = \gamma(\Delta z; \theta | \lambda) \prod_{i=1}^N \int_{\Delta z_i} \pi(\Delta\tilde{w}_i; \theta_i | z_i, \lambda) d\gamma(z_i; \theta_i | \lambda_i),$$

where

$$(2.3) \quad \pi(\Delta\tilde{w}; \theta | z, \lambda) = \int_B \delta_T(\Delta\tilde{w}; \theta | z, w, \lambda) d\kappa(w; \theta | z, \lambda),$$

and triples $(Z_1, W_1, \tilde{W}_1), \dots, (Z_N, W_N, \tilde{W}_N)$ and Z are assumed conditionally independent as before.

In the case of a system taught by estimating its actions Ξ represents the set of previously received inputs Z_1, \dots, Z_N , corresponding system outputs $\tilde{W}_1, \dots, \tilde{W}_N$, teacher outputs $\tilde{W}_1, \dots, \tilde{W}_N$, and the input Z to which the system must respond optimally. The system outputs $\tilde{W}_1, \dots, \tilde{W}_N$ represent random variables in B with conditional probability measures $\delta^1(\Delta\hat{w} | z), \dots, \delta^N(\Delta\hat{w} | z)$ independent of λ (since the values $\lambda_1, \dots, \lambda_N, \lambda$ of the random function $A(\theta)$ in the points $\theta_1, \dots, \theta_N, \theta$ remain

unknown to the system; otherwise learning would be of no sense). Yet each of the measures $\delta^2, \dots, \delta^N$ may depend on previously received teaching signals, namely $\delta^i(\Delta w | z)$ may depend also on the values $z_1, \dots, z_{i-1}, \hat{w}_1, \dots, \hat{w}_{i-1}, \tilde{w}_1, \dots, \tilde{w}_{i-1}$ of $Z_1, \dots, Z_{i-1}, \tilde{W}_1, \dots, \tilde{W}_{i-1}, \hat{W}_1, \dots, \hat{W}_{i-1}$ ($i = 2, \dots, N$). Hence the conditional probability measure of Ξ in the Cartesian product $X = A^{N+1} \times B^N \times \tilde{B}^N$ is given in this case by

$$(2.4) \quad \sigma(\Delta \xi | \lambda, \bar{\lambda}) = \gamma(\Delta z; \theta | \lambda) \prod_{i=1}^N \int_{\Delta z_i} d\gamma(z_i; \theta_i | \lambda_i) \times \\ \times \int_{\Delta \hat{w}_i} \pi(\Delta \hat{w}_i; \theta_i | z_i, \hat{w}_i, \lambda_i) d\delta^i(\hat{w}_i | z_i),$$

where

$$(2.5) \quad \pi(\Delta \hat{w}; \theta | z, \hat{w}, \lambda) = \int_B \delta_T(\Delta \hat{w}; \theta | z, w, \hat{w}, \lambda) d\kappa(w; \theta | z, \lambda),$$

and the quadruples $(Z_1, W_1, \hat{W}_1, \tilde{W}_1), \dots, (Z_N, W_N, \hat{W}_N, \tilde{W}_N)$ and Z are assumed conditionally independent, as before.

The system is required to elaborate the optimal estimate W^* of the desired output W at the first cycle after learning using the set Ξ of all signals received.

The average loss (i.e. the expected value of the loss function) at the first cycle after learning is determined by

$$(2.6) \quad R(\delta) = \mathbf{E}(W, \hat{W} | A) = \\ = \int_L d_\lambda \int_L d_{\bar{\lambda}} A(\lambda, \bar{\lambda}) \int_X d\sigma(\xi | \lambda, \bar{\lambda}) \int_B d\delta(\hat{w} | \xi) \int_B l(w, \hat{w} | \lambda) d\kappa(w; \theta | z, \lambda),$$

where $A(\Delta \lambda, \Delta \bar{\lambda})$ is the joint probability measure of the random variables $A(\theta), A(\theta_1), \dots, A(\theta_N)$ in the Cartesian product $L \times \bar{L} = L^{N+1}$, and $\delta(\Delta \hat{w} | \Delta)$ is the system decision function representing the conditional probability measure of system output \hat{W} at the first cycle after learning given the value ξ of Ξ .

The problem is then to find an optimal system decision function $\delta_{\text{op}}(\Delta \hat{w} | \xi)$ minimizing the average loss $R(\delta)$. The system using this optimal decision function will be the required *Bayes optimal learning system*.

To solve the problem we notice that for any sets $\Delta \xi, \Delta \lambda, \Delta \bar{\lambda}$ of values of $\xi, \lambda, \bar{\lambda}$ respectively

$$(2.7) \quad \int_{\Delta \lambda} d_\lambda \int_{\Delta \bar{\lambda}} d_{\bar{\lambda}} A(\lambda, \bar{\lambda}) \int_{\Delta \xi} d\sigma(\xi | \lambda, \bar{\lambda}) = \int_{\Delta \xi} d\beta(\xi) \int_{\Delta \lambda} d_\lambda \int_{\Delta \bar{\lambda}} d_{\bar{\lambda}} A(\lambda, \bar{\lambda} | \xi),$$

where

$$(2.8) \quad \beta(\Delta \xi) = \int_L d_\lambda \int_L d_{\bar{\lambda}} \sigma(\Delta \xi | \lambda, \bar{\lambda}) d_X A(\lambda, \bar{\lambda})$$

is the unconditional probability measure of Ξ , and

$$(2.9) \quad \eta(\Delta\lambda, \Delta\bar{\lambda} \mid \xi) = \int_{\Delta\lambda} d_\lambda \int_{\Delta\bar{\lambda}} \frac{d\sigma(\xi \mid \lambda, \bar{\lambda})}{d\beta(\xi)} d_{\bar{\lambda}} A(\lambda, \bar{\lambda})$$

the joint conditional probability measure of $A(\theta)$, $A(\theta_1)$, ..., $A(\theta_N)$ given the value ξ of Ξ . Using (2.7) and introducing the conditional probability measure of $A(\theta)$:

$$(2.10) \quad \Omega(\Delta\lambda \mid \xi) = \eta(\Delta\lambda, L \mid \xi) = \int_{\Delta\lambda} d_\lambda \int_L \frac{d\sigma(\xi \mid \lambda, \bar{\lambda})}{d\beta(\xi)} d_{\bar{\lambda}} A(\lambda, \bar{\lambda}),$$

(2.6) becomes

$$(2.11) \quad R(\delta) = \int_{\mathbf{X}} d\beta(\xi) \int_B \varrho(\xi, \hat{w}) d\delta(\hat{w} \mid \xi)$$

with

$$(2.12) \quad \varrho(\xi, \hat{w}) = \int_Z d\Omega(\lambda \mid \xi) \int_B l(w, \hat{w} \mid \lambda) dz(w; \theta \mid z, \lambda).$$

If for any ξ there exists a unique w^* satisfying

$$(2.13) \quad \varrho(\xi, w^*) \leq \varrho(\xi, \hat{w}) \quad \text{for all } \hat{w},$$

then the optimal system decision function $\delta_{\text{opt}}(\Delta\hat{w} \mid \xi)$ represents, obviously, the measure condensed in single point $w^* = Q\xi$ (Q being the deterministic operator establishing the correspondence between ξ and w^* satisfying (2.13)). Thus the Bayes optimal learning system is deterministic in this case and its output w^* corresponding to the value ξ of Ξ is determined as the unique value w^* of \hat{w} giving the least value to $\varrho(\xi, \hat{w})$.

If for any ξ there exists a set C_ξ of values w^* satisfying (2.13), then the optimal system decision function $\delta_{\text{opt}}(\Delta\hat{w} \mid \xi)$ represents an arbitrary probability measure on C_ξ . There exists then an infinite set of Bayes optimal learning systems, some of which are deterministic and others stochastic. Each of the latter includes a random mechanism of choice of w^* from C_ξ for any given value ξ of Ξ .

As we see from (2.12) the algorithm of learning of the Bayes optimal learning system consists of elaborating conditional probability measure $\Omega(\Delta\lambda \mid \xi)$ of $A(\theta)$ at the first cycle of system action after receiving the teaching signals. $\Omega(\Delta\lambda \mid \xi)$ represents in fact the *posterior* probability measure of $A(\theta)$ after receiving the teaching signals, while $A(\Delta\lambda, L)$ is the *prior* probability measure of $A(\theta)$.

It should be emphasized that any Bayes optimal learning system is optimal for any value ξ of Ξ , i.e. for any given set of teaching signals and any input Z after receiving the teaching signals, the expected loss $\varrho(\xi, w^*)$ having the least possible value for any value of ξ .

Let us now estimate the effect of learning. To do this we compare the Bayes optimal learning system with the corresponding Bayes optimal system with complete information about the parameter λ . This Bayes optimal system minimizes the conditional loss

$$(2.14) \quad \varrho_\lambda(z, \hat{w}) = \int_B l(w, \hat{w} \mid \lambda) d\mathcal{K}(w; \theta \mid z, \lambda).$$

If for any z there exists such w_λ^* for which

$$(2.15) \quad \varrho_\lambda(z, w_\lambda^*) \leq \varrho_\lambda(z, \hat{w}) \quad \text{for all } \hat{w},$$

then the optimal decision function $\delta_\lambda(\Delta \hat{w} \mid z)$ is condensed in single point $\hat{w} = w_\lambda^* = A(\lambda)z$, $A(\lambda)$ being an operator depending on λ (and θ , certainly). If for any z there exists a set $C_{z,\lambda}$ of w_λ^* satisfying (2.15), then the optimal decision function $\delta_\lambda(\Delta \hat{w} \mid z)$ represents an arbitrary probability measure on $C_{z,\lambda}$.

In both cases (2.15) is satisfied for $\hat{w} = w^* = Q\xi$. Hence the relative amount of the loss due to the absence of complete information about the value of the parameter λ :

$$(2.16) \quad \varepsilon_{\lambda\xi}(\delta_{\text{opt}}) = \frac{\varrho_\lambda(z, w^*) - \varrho_\lambda(z, w_\lambda^*)}{\varrho_\lambda(z, w_\lambda^*)}$$

can be taken as the measure of closeness of the Bayes optimal learning system to the Bayes optimal system with complete information.

To estimate the average performance of the learning system for all possible values of ξ given λ , we derive from (2.15)

$$(2.17) \quad \int_X d\sigma(\xi \mid \lambda, \bar{\lambda}) \int_B \varrho_\lambda(z, \hat{w}) d\delta_\lambda(\hat{w} \mid z) \leq \int_X d\sigma(\xi \mid \lambda, \bar{\lambda}) \int_B \varrho_\lambda(z, \hat{w}) d\delta_{\text{opt}}(\hat{w} \mid \xi).$$

Averaging with respect to $\bar{\lambda}$ and denoting the conditional average loss corresponding to a given value of the parameter λ as

$$(2.18) \quad r_\lambda(\delta) = \int_L d_x \frac{d_\lambda A(\lambda, \bar{\lambda})}{d_\lambda \int_L d_{\bar{\gamma}} A(\lambda, \bar{\gamma})} \int_X d\sigma(\xi \mid \lambda, \bar{\lambda}) \int_B \varrho_\lambda(z, \hat{w}) d\delta(\hat{w} \mid \xi),$$

we obtain from (2.17) $r_\lambda(\delta_\lambda) \leq r_\lambda(\delta_{\text{opt}})$. So the relative amount of the conditional average loss

$$(2.19) \quad \varepsilon_\lambda(\delta_{\text{opt}}) = \frac{r_\lambda(\delta_{\text{opt}}) - r_\lambda(\delta_\lambda)}{r_\lambda(\delta_\lambda)}$$

can serve as a measure of performances of the Bayes optimal learning system for a given value of λ .

Finally, averaging (2.17) with respect to both $\bar{\lambda}$ and λ , we obtain $R(\delta_{\lambda}) \leq R(\delta_{opt})$. So the relative amount of the average loss

$$(2.20) \quad \varepsilon(\delta_{opt}) = \frac{R(\delta_{opt}) - R(\delta_{\lambda})}{R(\delta_{\lambda})}$$

can be taken as a total measure of performance of the Bayes optimal learning system (the measure of its *potential learnability*) [2].*

Similarly the quantities $\varepsilon_{\lambda\xi}(\delta)$, $\varepsilon_{\lambda}(\delta)$ and $\varepsilon(\delta)$ for any $\delta(\Delta\tilde{w} | \xi)$ can serve as respective measures of the result of learning for any algorithm of learning.

To estimate the result of learning it is sufficient to compare the values of $\varepsilon_{\lambda\xi}$, ε_{λ} or ε for the Bayes optimal learning system with the respective values of $\varepsilon_{\lambda\xi}$, ε_{λ} , ε for the Bayes optimal non-learning system which has no information about λ and does not receive teaching signals. To find the decision function $\delta'_{opt}(\Delta\tilde{w} | z)$ of such a system it is sufficient to replace $\sigma(\Delta\xi | \lambda, \bar{\lambda})$ in all previous formulas by $\gamma(\Delta z; \theta | \lambda)$.

3. SOME GENERAL PROPERTIES OF OPTIMAL LEARNING PROCESSES

From (2.12), (2.10), (2.1), (2.2), (2.4) follow some general properties of optimal learning processes. Primarily, the Bayes optimal learning system represents a permanently learning system, namely it never ceases self-learning, as shown by (2.1), (2.2), (2.4) and (2.10).

It is clear that the learning of a system is possible only if the random variables $A(\theta_1), \dots, A(\theta_N), A(\theta)$ are interdependent, i.e. if the unknown parameter λ varies sufficiently slowly. The learning is impossible, if $A(\theta_1), \dots, A(\theta_N), A(\theta)$ are statistically independent, since $\Omega(\Delta\lambda | \xi)$ depends only on the value z of the input Z , and is quite independent of previously obtained teaching signals in such a case. Closer is the interdependence between $A(\theta_1), \dots, A(\theta_N), A(\theta)$ to some determined functional relation, stronger is the effect of previously obtained teaching signals on $\Omega(\Delta\lambda | \xi)$, and therefore better are the results of learning. The most effective learning is attained when $A(\theta_1) = \dots = A(\theta_N) = A(\theta)$, i.e. when the unknown parameter λ remains constant during learning and subsequent system action.

In the special case of the ideal teacher $\delta_T(\Delta\tilde{w}; \theta | z, w, \lambda)$ is unity for any set $\Delta\tilde{w}$ including w , and zero otherwise. So $\pi(\Delta\tilde{w}; \theta | z, \lambda) = \varkappa(\Delta\tilde{w}; \theta | z, \lambda)$ in this case, and (2.2) becomes

$$(3.1) \quad \sigma(\Delta\xi | \lambda, \bar{\lambda}) = \gamma(\Delta z; \theta | \lambda) \prod_{i=1}^N \int_{\Delta z_i} \varkappa(\Delta\tilde{w}_i; \theta_i | z_i, \lambda_i) d\gamma(z_i; \theta_i | \lambda_i).$$

* If $Q_2(z, w^*) = 0$ for some z, λ , then the loss $Q_2(z, w^*)$ should be taken as a measure of performance of the learning system for these z, λ instead of $\varepsilon_{\lambda\xi}(\delta_{opt})$. Similarly, if $r_{\lambda}(\delta_{\lambda}) = 0$ for some λ , then the conditional average loss $r_{\lambda}(\delta_{opt})$ can serve as a measure of performance instead of $\varepsilon_{\lambda}(\delta_{opt})$. Finally, if $R(\delta_{\lambda}) = 0$, the average loss $R(\delta_{opt})$ would be a suitable measure of performance instead of $\varepsilon(\delta_{opt})$.

In the case of a real teacher δ_T in (2.3) and (2.5) is independent of w , and therefore $\pi = \delta_T$, and (2.2) and (2.4) take respectively the forms

$$(3.2) \quad \sigma(\Delta\xi | \lambda, \lambda) = \gamma(\Delta z; \theta | \lambda) \prod_{i=1}^N \int_{\Delta z_i} \delta_T(\Delta\tilde{w}_i; \theta_i | z_i, \lambda_i) d\gamma(z_i; \theta_i | \lambda_i),$$

$$(3.3) \quad \sigma(\Delta\xi | \lambda, \lambda) = \gamma(\Delta z; \theta | \lambda) \prod_{i=1}^N \int_{\Delta z_i} d\gamma(z_i; \theta_i | \lambda_i) \times \\ \times \int_{\Delta\tilde{w}_i} \delta_T(\Delta\tilde{w}_i; \theta_i | z_i, \tilde{w}_i, \lambda_i) d\delta^i(\tilde{w}_i | z_i).$$

The comparison of (3.1) with (3.2) shows that the real teacher teaching the system by show with $\delta_T(\Delta\tilde{w}; \theta | z, \lambda) = \varkappa(\Delta\tilde{w}; \theta | z, \lambda)$ is completely equivalent to the ideal teacher. In other words the real teacher whose decision function coincides with the conditional distribution of the desired output W for a given value z of the input Z , is completely equivalent to the ideal teacher from the point of view of its teaching capabilities.

If a real teacher teaching the system by show represents a deterministic system, then the measure δ_T is condensed in single point $\tilde{w} = A_T(\lambda)z$, $A_T(\lambda)$ being some deterministic operator depending on the corresponding value λ of $A(\theta)$ (and may be on θ also), then the posterior measure $\eta(\Delta\lambda, \Delta\tilde{\lambda} | \xi)$ differs from zero only on the subset of $L \times \bar{L}$ determined by

$$(3.4) \quad \tilde{w}_i = A_T(\lambda_i) z_i \quad (i = 1, \dots, N).$$

If the random variables $A(\theta_1), \dots, A(\theta_N), A(\theta)$ are strongly correlated, then (3.4) imposes rather strong restriction on the domain of possible values of $A(\theta)$ owing to which the a posteriori measure $\Omega(\Delta\lambda | \xi)$ will be condensed in a narrow domain near the unknown true value of λ . The effect of learning may be better in this case than in the case of the ideal teacher.

If, in particular $A(\theta_1) = \dots = A(\theta_N) = A(\theta)$, i.e. the parameter λ is constant, then the posterior measure $\Omega(\Delta\lambda | \xi)$ differs from zero only on the subset of L determined by the equations

$$(3.5) \quad \tilde{w}_i = A_T(\lambda) z_i \quad (i = 1, \dots, N).$$

If there exists a finite number, say r , of such equations having a unique solution with respect to λ , then the unknown parameter λ is exactly determined after receiving r pairs $(z_1, \tilde{w}_1), \dots, (z_r, \tilde{w}_r)$ of a teaching signals, and the optimal learning system becomes the Bayes optimal system with complete information about λ , for which $\varepsilon_{\lambda\xi}(\delta_{opt}) = 0$.

Thus the deterministic real teacher whose operator $A_T(\lambda)$ admits a finite number of equations of the form (3.5) having a unique solution with respect to λ , is the best

possible teacher making the system completely learned after showing to it the respective finite number of teaching pairs of signals. The ideal teacher is thus by no means the best one in such a case. This result seems to be a paradox: the ideal teacher showing to the system *exact* values of the desired output corresponding to given values of the input is worse than the real teacher *committing errors* with probability one. The explanation of this paradox is very simple: since in the general case the input Z represents some signal distorted by noises, and the desired output W depends only on the signal, there exists no deterministic operator transforming Z into W . This is the reason why the relation between Z and W is much more difficult to determine from observations than the relation between Z and the output \tilde{W} of a deterministic teacher.

It is also clear that in such cases stochastic real teachers with sufficiently small variances of \tilde{W} can be better than the ideal teacher.

In special cases where the noises are absent or have small variances, the ideal teacher may certainly be better than any real teacher. The examples where the ideal teacher can make a system completely learned after showing to it a finite number of pairs input-output are numerous in pattern recognition theory (see, e.g., [29]).

Let us now consider an important special case where $A(\theta)$ represents a Markov random process. Let us partition the teaching cycles and the corresponding values $\theta_1, \dots, \theta_N$ of θ into two groups $\theta_1, \dots, \theta_K$ and $\theta_{K+1}, \dots, \theta_N$ and let Ξ_1 be the set of teaching signals received in the first K cycles plus the $(K + 1)$ -th input Z_{K+1} , Ξ_2 the set of all remaining teaching signals plus the input Z at the first cycle after learning, λ_i^j the set $\lambda_i, \dots, \lambda_j$, L_i^j the Cartesian product L^{j-i+1} , $\sigma_1(\Delta\xi_1 | \lambda_{K+1}, \lambda_1^K)$, $\sigma_2(\Delta\xi_2 | \lambda, \lambda_{K+1}^N)$, $\beta_2(\Delta\xi_2)$ the respective probability measures σ and β . Then taking into account that in this case

$$(3.6) \quad A(\Delta\lambda, \Delta\tilde{\lambda}) = \int_{\Delta\lambda_{K+1}} A_2(\Delta\lambda, \Delta\lambda_{K+2}^N | \lambda_{K+1}) d_{\lambda_{K+1}} A_1(\lambda_{K+1}, \Delta\lambda_1^K),$$

$A_1(\Delta\lambda_{K+1}, \Delta\lambda_1^K)$ being the probability measure of $A(\theta_1), \dots, A(\theta_{K+1})$, and $A_2(\Delta\lambda, \Delta\lambda_{K+2}^N | \lambda_{K+1})$ the conditional probability measure of $A(\theta_{K+2}), \dots, A(\theta_N)$, $A(\theta)$ given the value λ_{K+1} of $A(\theta_{K+1})$, we can rewrite (2.10) as

$$(3.7) \quad \Omega(\Delta\lambda | \xi) = \int_{\Delta\lambda} d_\lambda \int_{L^{N-K+1}} \frac{d\sigma_2(\xi_2 | \lambda, \lambda_{K+1}^N)}{d\beta_2(\xi_2)} d_{\lambda_{K+1}} A_3(\lambda, \lambda_{K+1}^N)$$

with

$$(3.8) \quad A_3(\Delta\lambda, \Delta\lambda_{K+1}^N) = \int_{\Delta\lambda_{K+1}} A_2(\Delta\lambda, \Delta\lambda_{K+2}^N | \lambda_{K+1}) d_{\lambda_{K+1}} \Omega_1(\lambda_{K+1} | \xi_1),$$

$$\Omega_1(\Delta\lambda_{K+1} | \xi_1) = \int_{\Delta\lambda_{K+1}} d_{\lambda_{K+1}} \int_{L^{1,K}} \frac{d\sigma_1(\xi_1 | \lambda_{K+1}, \lambda_1^K)}{d\beta_1(\xi_1)} d_{\lambda_{1,K}} A_1(\lambda_{K+1}, \lambda_1^K).$$

Formulas (3.7) and (3.8) show that the learning process can be realized recursively in this case, namely $\Omega(\Delta\lambda | \xi)$ can be calculated step by step after each teaching cycle using the posterior distribution previously obtained as the prior distribution. This fact permits to reduce essentially the capacity of the memory of computers.

Of course, this is also true for the special case of constant unknown parameter λ .

All the properties of the optimal learning processes studied above were previously established for the special case of constant unknown parameter λ in [1-3].

We have restricted ourselves to studying learning processes for any finite number N of teaching cycles. The reason of this is that only a finite period of learning can be realized in practice. Yet the questions of convergence of learning processes when $N \rightarrow \infty$ are also of interest from the theoretical point of view. It is clear from the definitions of Section 2 that the optimal learning process is convergent for a given value of λ and a given sequence ξ of teaching signals, if $\varepsilon_{\lambda\xi}(\delta_{\text{opt}}) \rightarrow 0$ when $N \rightarrow \infty$. It is convergent for a given value of λ and almost all possible teaching sequences ξ , if $\varepsilon_{\lambda}(\delta_{\text{opt}}) \rightarrow 0$. Finally, it is convergent for almost all values of λ and almost possible sequences ξ , if $\varepsilon(\delta_{\text{opt}}) \rightarrow 0$.

4. CASE OF DISCRETE SYSTEMS

All the results obtained in previous sections are valid for the most general forms of inputs and outputs of systems. They may be elements of arbitrary sets. The cost of this generality is practical impossibility to use the theory for direct calculations. Yet the input and output of any real technical system represent some ordinary scalar or vector functions of time or other arguments. The large class of modern technical systems is the class of discrete or sampled-data systems. This class include, for instance, all digital computers. The inputs and outputs of such systems represent functions of discrete arguments taking always only finite number of values. Hence the sets of all the values of inputs and outputs may be considered as finite-dimensional vectors, and the sets A and B of the preceding theory represent ordinary finite-dimensional Euclidean spaces.

Thus in the case of discrete systems we may consider the input Z , desired and actual outputs of a system W , \tilde{W} , and the teacher output \tilde{W} as finite-dimensional random vectors and define their probability measures $\gamma(\Delta z; \theta | \lambda)$, $\varkappa(\Delta w; \theta | z, \lambda)$, $\delta_T(\Delta \tilde{w}; \theta | z, w, \tilde{w}, \lambda)$ by respective probability densities*:

$$(4.1) \quad \gamma(\Delta z; \theta | \lambda) = \int_{\Delta z} g(z; \theta | \lambda) dz,$$

* If Z , W and \tilde{W} represent discrete random variables, then the probability densities g , k and d_T represent linear combinations of Dirac delta-functions. If Z , W and \tilde{W} have some possible values with non-zero probabilities besides continuous domains of possible values, then the respective probability densities g , k and d_T represent sums of functions integrable in Riemann sense and linear combinations of delta-functions.

$$(4.2) \quad \alpha(\Delta w; \theta | z, \lambda) = \int_{\Delta w} k(w; \theta | z, \lambda) dw,$$

$$(4.3) \quad \delta_T(\Delta \tilde{w}; \theta | z, w, \hat{w}, \lambda) = \int_{\Delta \tilde{w}} d_T(\tilde{w}; \theta | z, w, \hat{w}, \lambda) d\tilde{w}.$$

To obtain results applicable for direct calculations in technical system design we suppose that the parameter λ represents also a finite-dimensional vector, and, consequently, L is an ordinary Euclidean space. Thus $A(\theta)$ of the preceding theory represents now a random vector function of θ , and the probability measures $A(\Delta\lambda, \Delta\tilde{\lambda})$ and $\Omega(\hat{w}\lambda | \xi)$ can be defined by respective probability densities:

$$(4.4) \quad A(\Delta\lambda, \Delta\tilde{\lambda}) = \int_{\Delta\lambda} \int_{\Delta\tilde{\lambda}} \alpha(\lambda, \tilde{\lambda}) d\tilde{\lambda},$$

$$(4.5) \quad \Omega(\Delta\lambda | \xi) = \int_{\Delta\lambda} \omega(\lambda | \xi) d\lambda.$$

Using (4.1) – (4.5) we obtain from (2.10), (2.1), (2.2) and (2.4)

$$(4.6) \quad \omega(\lambda | \xi) = c(\xi) g(z; \theta | \lambda) \int_L \alpha(\lambda, \tilde{\lambda}) \prod_{i=1}^N g(z_i; \theta_i | \lambda_i) \times \\ \times p(\tilde{w}_i; \theta_i | z_i, \hat{w}_i, \lambda_i) d\tilde{\lambda},$$

where integration extends over the domain of all possible values of the composed random vector $\tilde{\lambda}$ formed by all the components of the vectors $A(\theta_1), \dots, A(\theta_N)$, the function $p(\tilde{w}; \theta | z, \hat{w}, \lambda)$ is unity in the case of self-learning, and is determined by

$$(4.7) \quad p(\tilde{w}; \theta | z, \hat{w}, \lambda) = \int_B d_T(\tilde{w}; \theta | z, w, \hat{w}, \lambda) k(w; \theta | z, \lambda) dw$$

in the case of learning with a teacher, and

$$(4.8) \quad c(\xi) = \left[\int_L g(z; \theta | \lambda) d\lambda \int_L \alpha(\lambda, \tilde{\lambda}) \prod_{i=1}^N g(z_i; \theta_i | \lambda_i) \times \right. \\ \left. \times p(\tilde{w}_i; \theta_i | z_i, \hat{w}_i, \lambda_i) d\tilde{\lambda} \right]^{-1}$$

represents normalizing constant generally depending on ξ .

In the special case of the ideal teacher, d_T represents Dirac delta-function (certainly, multidimensional, in general), and $p = k$. In the case of any real teacher d_T is independent of w , and $p = d_T$.

The value w^* of the output of the Bayes optimal learning system corresponding to a given value z of the input Z is determined in this case as the value of \hat{w} minimizing

$$(4.9) \quad \varrho(\xi, \hat{w}) = \int_L \omega(\lambda | \xi) d\lambda \int_B l(w, \hat{w} | \lambda) k(w; \theta | z, \lambda) dw,$$

whereas the value w_λ^* of the output of the Bayes optimal system with complete information about λ corresponding to the same value z of the input Z is determined as the value of \hat{w} minimizing

$$(4.10) \quad \varrho_\lambda(z, \hat{w}) = \int_B l(w, \hat{w} | \lambda) k(w; \theta | z, \lambda) dw.$$

The formulas (2.18) and (2.6) determining the conditional average loss $r_\lambda(\delta)$ and average loss $R(\delta)$ become in this case respectively

$$(4.11) \quad r_\lambda(\delta) = \int_L \alpha_2(\bar{\lambda} | \lambda) d\bar{\lambda} \int_X s(\xi | \lambda, \bar{\lambda}) d\xi \int_B \varrho_\lambda(z, \hat{w}) d(\hat{w} | \xi) d\hat{w},$$

$$(4.12) \quad R(\delta) = \int_L r_\lambda(\delta) \alpha_1(\lambda) d\lambda,$$

where

$$(4.13) \quad \alpha_1(\lambda) = \int_L \alpha(\lambda, \bar{\lambda}) d\bar{\lambda}, \quad \alpha_2(\bar{\lambda} | \lambda) = \frac{\alpha(\lambda, \bar{\lambda})}{\alpha_1(\lambda)},$$

$$(4.14) \quad s(\xi | \lambda, \bar{\lambda}) = g(z; \theta | \lambda) \prod_{i=1}^N g(z_i; \theta_i | \lambda_i) \times d^i(\hat{w}_i | z_i) p(\bar{w}_i; \theta_i | z_i, \hat{w}_i, \lambda_i),$$

and $d^i(\hat{w} | z), \dots, d^N(\hat{w} | z), d(\hat{w} | \xi)$ represent the conditional probability densities of system output \hat{W} at the respective cycles of its action determining its decision functions, of which $d(\hat{w} | \xi)$ is to be optimized.

Let us now consider the special case of learning pattern recognition systems designed to recognize which of the mutually exclusive patterns A_1, \dots, A_n is present in the input Z received.* The desired output W represents in this case the number of the pattern A_w , and therefore is a discrete random variable with n possible values $1, \dots, n$. The actual output of the system \hat{W} represents the number of the pattern determined by the system. If the refusal to decide is admissible, then \hat{W} is a random variable with $n + 1$ possible values $0, 1, \dots, n$, the value 0 being assigned to the refusal. As to the teacher output \bar{W} , it may represent either the number of the pattern determined by the teacher, or the number of the group of patterns, if the teacher shows to the system only to which of several groups into which the patterns are partitioned the input Z corresponds, or the estimate of the system response given by the teacher in the case of teaching by estimating system actions. Supposing in the latter case that the teacher estimate can assume only integer values $0, 1, \dots, r$, we cover all the cases assuming that the teacher output \bar{W} represents a random variable with integer possible values $0, 1, \dots, r$ (r being equal to n in the first case).

* We call *pattern* any set of subjects having some common features in virtue of which they are considered as belonging to the same class. The pattern recognition problem consists in distinguishing between the subjects of different classes, i.e. of classifying objects into several classes.

It should be emphasized that this case covers also the problem of recognizing any finite sequence of patterns, since any such sequence may be considered as a single composed pattern.

Let p_1, \dots, p_n be the probabilities of the patterns A_1, \dots, A_n respectively, $f_w(z; \theta | \lambda)$ the conditional probability density of the input Z for a given pattern A_w . Then

$$(4.15) \quad g(z; \theta | \lambda) = \sum_{k=1}^n p_k f_k(z; \theta | \lambda),$$

$$(4.16) \quad k(w; \theta | z, \lambda) = \sum_{k=1}^n P_k(z, \lambda, \theta) \delta(w - k),$$

$$(4.17) \quad d_T(\tilde{w}; \theta | z, w, \hat{w}, \lambda) = \sum_{h=1}^r Q_h(z, w, \hat{w}, \lambda, \theta) \delta(\tilde{w} - h),$$

and (4.7) becomes

$$(4.18) \quad p(\tilde{w}; \theta | z, \hat{w}, \lambda) = \sum_{h=0}^r \Psi_h(z, \hat{w}, \lambda, \theta) \delta(\tilde{w} - h),$$

where

$$(4.19) \quad \Psi_h(z, \hat{w}, \lambda, \theta) = \sum_{k=1}^n P_k(z, \lambda, \theta) Q_h(z, k, \hat{w}, \lambda, \theta).$$

Formula (4.6) takes the form

$$(4.20) \quad \omega(\lambda | \xi) = c(\xi) g(z; \theta | \lambda) \int_L \alpha(\lambda, \bar{\lambda}) \prod_{i=1}^N g(z_i; \theta_i | \lambda_i) \Psi_{\tilde{w}_i}(z_i, \hat{w}_i, \lambda_i, \theta_i) d\bar{\lambda},$$

where

$$(4.21) \quad c(\xi) = \left[\int_L g(z; \theta | \lambda) d\lambda \int_L \alpha(\lambda, \bar{\lambda}) \prod_{i=1}^N g(z_i; \theta_i | \lambda_i) \Psi_{\tilde{w}_i}(z_i, \hat{w}_i, \lambda_i, \theta_i) d\bar{\lambda} \right]^{-1}.$$

To obtain this formula from (4.6) strictly the function $\delta(x)$ should be replaced by some bounded function different from zero only in a narrow interval $(-\varepsilon, \varepsilon)$, $\varepsilon < 1/2$. Evidently this does not alter the problem. The this function will be cancelled in (4.20) and (4.21), and we obtain (4.20) in the limit when $\varepsilon \rightarrow 0$. Formula (4.20) can, of course, be obtained directly from the general formulas (2.10), (2.1), (2.2) and (2.4) without limiting processes.

Finally, (4.9) takes the form

$$(4.22) \quad \varrho(\xi, h) = \sum_{k=1}^n I_{kh} P_k^*(\xi, \lambda) \quad (h = \hat{w} = 0, 1, \dots, n),$$

where

$$(4.23) \quad P_k^*(\xi, \lambda) = \int_L P_k(z, \lambda, \theta) \omega(\lambda | \xi) d\lambda \quad (k = 1, \dots, n).$$

362 The Bayes optimal learning system must choose the value h of \hat{w} minimizing $q(\zeta, h)$, whereas the Bayes optimal system with complete information about λ must choose the value m minimizing

$$(4.24) \quad \varrho_\lambda(z, m) = \sum_{k=1}^n l_{km} P_k(z, \lambda, \theta) \quad (m = 0, 1, \dots, n).$$

Example. Let us consider a learning system designed to estimate the parameter U of the input

$$(4.25) \quad Z(t) = U \varphi(t) + X(t),$$

where $\varphi(t)$ is a given function, U normally distributed random variable with unknown expected value λ and known variance D_u , $X(t)$ normally distributed random function independent of U with zero expected value and known covariance $K_x(t, t')$. The performance of the system is measured by its mean square error $E\{\hat{W} - W\}^2$.

Considering the case of discrete systems we assume that the input Z excites the system at time instants $t = t_1, t_2, \dots, t_n$, and the optimal estimate of U is required at the same time instants. The desired output W is in this case the parameter U . The Bayes optimal system with complete information about λ elaborates its output determined by

$$(4.26) \quad W_\lambda^*(t_v) = \frac{D_u H^{(v)} + \lambda}{b^{(v)} D_u + 1} \quad (v = 1, \dots, n),$$

where

$$(4.27) \quad H^{(v)} = \sum_{\mu=1}^v g_{v\mu} Z(t_\mu), \quad b^{(v)} = \sum_{\mu=1}^v g_{v\mu} \varphi(t_\mu),$$

the coefficients $g_{v\mu}$ being determined by the sets of linear algebraic equations

$$(4.28) \quad \sum_{\mu=1}^v g_{v\mu} K_x(t_\mu, t_\sigma) = \varphi(t_\sigma) \quad (\sigma = 1, \dots, v; v = 1, \dots, n).$$

To determine the Bayes optimal learning system we suppose that the system receives at N cycles of the period of learning the values z_1, \dots, z_N of the input and the respective values $\tilde{w}_1, \dots, \tilde{w}_N$ of teacher output, and then must elaborate the optimal estimate of the signal parameter U at every instant t_1, \dots, t_n of the first cycle after learning. As to the teacher we shall suppose that it elaborates at each cycle the Bayes optimal estimates $W_\lambda^*(t_v)$ ($v = 1, \dots, n$) of $W = U$ with random errors $\tilde{Y}(t_1), \dots, \tilde{Y}(t_n)$ which represent the values of a normally distributed random function $\tilde{Y}(s)$ with zero expected value and known covariance $K_Y(s, s')$.

The probability density of the input Z and the decision function of the teacher at each of N teaching cycles are determined in this case by

$$(4.29) \quad g(z_i; \varepsilon_i | \lambda) = (2\pi)^{-n/2} |K_n|^{-1/2} \times \\ \times \exp \left\{ -\frac{1}{2} [z_i - \zeta(\lambda)]' K_n^{-1} [z_i - \zeta(\lambda)] \right\},$$

$$(4.30) \quad d_T(\tilde{w}_i; \varepsilon_i | z_i, \lambda) = (2\pi)^{-n/2} |K_T|^{-1/2} \times \\ \times \exp \left\{ -\frac{1}{2} [\tilde{w}_i - v_i(\lambda)]' K_T^{-1} [\tilde{w}_i - v_i(\lambda)] \right\},$$

where z_i represents $n \times 1$ matrix whose elements are the values $z_i(t_v)$ ($v = 1, \dots, n$) of the input $Z_i(t_v)$ at the i -th cycle, $\zeta(\lambda)$ the $n \times 1$ matrix with the elements $\zeta^{(v)}(\lambda) = \lambda \varphi(t_v)$ ($v = 1, \dots, n$), K_n the $n \times n$ covariance matrix of the input Z with the elements $D_u \varphi(t_v) \varphi(t_\mu) + K_x(t_v, t_\mu)$, \tilde{w}_i the $n \times 1$ matrix whose elements are the values of teacher output

$$(4.31) \quad \tilde{W}(t_v) = \frac{D_u H^{(v)} + \lambda}{b^{(v)} D_u + 1} + \tilde{Y}(t_v) \quad (v = 1, \dots, n)$$

at the i -th cycle, $v_i(\lambda)$ the $n \times 1$ matrix with the elements

$$(4.32) \quad v_i^{(v)}(\lambda) = \frac{D_u \eta_i^{(v)} + \lambda}{b^{(v)} D_u + 1}, \quad \eta_i^{(v)} = \sum_{\mu=1}^v g_{v\mu} z_i(t_\mu),$$

K_T the conditional covariance matrix of the random vector \tilde{W}_i with the elements $K_T(t_v, t_\mu)$. The probability density of the input Z at each instant t_v of the first cycle after learning is given by formula (4.29) with n replaced by v ($v = 1, \dots, n$).

Assuming a normal prior distribution of the unknown parameter considered as a random variable A with the expected value m and variance D , (4.6) and (4.7) yield normal posterior distribution for A , the posterior expected value $\lambda^*(t_v)$ and variance $\Delta(t_v)$ of A being determined by

$$(4.33) \quad \lambda^*(t_v) = A(t_v) \left\{ \frac{\eta^{(v)}}{b^{(v)} D_u + 1} + \frac{1}{b^{(n)} D_u + 1} \sum_{i=1}^N \eta^{(n)} + \sum_{i=1}^N [\zeta_i - \Phi(\eta_i)] + \frac{m}{D} \right\},$$

$$(4.34) \quad A(t_v) = \frac{D(b^{(v)} D_u + 1)(b^{(n)} D_u + 1)}{(b^{(v)} D_u + 1)[(1 + N\kappa D)(b^{(n)} D_u + 1) + Nb^{(n)} D] + b^{(v)} D(b^{(n)} D_u + 1)},$$

where, in addition to former notations,

$$(4.35) \quad \eta^{(v)} = \sum_{\mu=1}^v g_{v\mu} z(t_\mu), \quad \zeta_i = \sum_{v=1}^n h_v \tilde{w}_i(t_v),$$

$$(4.36) \quad \Phi(\eta) = \sum_{v=1}^n \frac{h_v \eta^{(v)}}{b^{(v)} D_u + 1}, \quad \kappa = \sum_{v=1}^n \frac{h_v}{b^{(v)} D_u + 1},$$

$z(t_1), \dots, z(t_n)$ represent the values of the input $Z(t)$ at the first cycle after learning, and the coefficients h_1, \dots, h_n are determined by the set of linear algebraic equations

$$(4.37) \quad \sum_{v=1}^n h_v K_T(t_v, t_\sigma) = \frac{1}{b^{(\sigma)} D_u + 1} \quad (\sigma = 1, \dots, n).$$

The output of the Bayes optimal learning system is given by

$$(4.38) \quad W^*(t_v) = \frac{D_u H^{(v)} + \lambda^*(t_v)}{b^{(v)} D_u + 1}.$$

364 Formula (2.20) gives in this case the following expression for the relative amount of the average loss at each step of the first cycle after learning:

$$(4.39) \quad \varepsilon(\delta_{\text{opt}}) = \frac{A(t_v)}{D_u(b^{(v)}D_u + 1)} \quad (v = 1, \dots, n).$$

For the Bayes optimal non-learning system $\lambda^*(t_v)$ and $A(t_v)$ must be replaced by the respective prior values m and D . Formula (4.39) shows then that the performance of the learning system exceeds that of the non-learning system as D exceeds $A(t_v)$.

As (4.34) shows $A(t_v) \rightarrow 0$ when $N \rightarrow \infty$, $\lambda^*(t_v)$ tends to the unknown true value of λ , and $\varepsilon(\delta_{\text{opt}}) \rightarrow 0$. The optimal learning process is thus convergent in this case.

The case of self-learning may be obtained as a special case where the variance of the teacher output W is infinite. Equations (4.37) give in this case $h_1 = \dots = h_n = 0$, and hence $\zeta_i = \Phi(\eta_i) = \alpha = 0$.

The case of the ideal teacher, as we proved, is formally the same as the case of the real teacher whose decision function d_T coincides with the conditional probability density of the desired output W given the input Z . Thus to obtain the case of the ideal teacher, we must put

$$(4.40) \quad K_{\mathcal{F}}(t_v, t_\mu) = \frac{D_u}{b^{(\sigma)}D_u + 1}, \quad \sigma = \max\{v, \mu\}.$$

Finally the case of the deterministic teacher which represents the Bayes optimal system with complete information about λ is obtained by putting $K_{\mathcal{F}}(s, s') = 0$. Equations (4.37) give in this case $h_1 = \dots = h_n = \infty$ after which (4.36) yields $\alpha = \infty$, and formula (4.34) shows that $A(t_v) = 0$ for any N and v in this case. This means that the system is completely learned and becomes itself the Bayes optimal system with complete information about λ after receiving a single value (at any of the points t_1, \dots, t_n) of the teacher output. This fact is quite clear, since $\tilde{Y}(s) = 0$ in this case, and (4.31) gives

$$(4.41) \quad \tilde{w}(t_v) = \frac{D_u \eta^{(v)} + \lambda}{b^{(v)}D_u + 1} \quad (v = 1, \dots, n).$$

Any of these equations, say the first one, yields immediately the exact value of the unknown parameter λ . This example illustrates very well the fact established above that the ideal teacher showing to the system exact value of the desired output, i.e. the exact value of the parameter U to be estimated, is worse than the deterministic real teacher estimating U with random errors.

5. CASE OF CONTINUOUS SYSTEMS

In the case of systems with continuous input Z represents a random function of the argument t continuously varying in a certain domain T . The outputs W , \hat{W} and \tilde{W} may be discrete in this case, i.e. random vectors as in the preceding section. The random function $Z(t)$ is scalar in the case of systems with one input, and vector in the case of systems with several inputs. The argument t represents usually time, but may be any scalar or vector variable as well.

To obtain practically applicable results we shall suppose as before that λ is a finite-

dimensional vector. Then from (2.10), (2.1)–(2.5), (4.2)–(4.5) and (4.7) follows

$$(5.1) \quad \omega(\lambda \mid \xi) = \int_L \frac{d\sigma(\xi \mid \lambda, \bar{\lambda})}{d\beta(\xi)} \alpha(\lambda, \bar{\lambda}) d\bar{\lambda}$$

with

$$(5.2) \quad \frac{d\sigma(\xi \mid \lambda, \bar{\lambda})}{d\beta(\xi)} = \frac{\prod_{i=1}^N p(\bar{w}_i; \theta_i \mid z_i, \hat{w}_i, \lambda_i)}{\int_L d\mu \int_L \frac{d\gamma(z; \theta \mid \mu)}{d\gamma(z; \theta \mid \lambda)} \prod_{i=1}^N \frac{d\gamma(z_i; \theta_i \mid \mu_i)}{d\gamma(z_i; \theta_i \mid \lambda_i)} p(\bar{w}_i; \theta_i \mid z_i, \hat{w}_i, \mu_i) \alpha(\mu, \bar{\mu}) d\bar{\mu}}$$

Thus it is sufficient to find Radon-Nikodym derivative $d\gamma(z; \theta \mid \mu)/d\gamma(z; \theta \mid \lambda)$.

We shall consider an important special case where $Z(t)$ statistically depends on a finite-dimensional random vector U and has normal conditional distribution for any possible values λ, u of A, U . In this case

$$(5.3) \quad \gamma(\hat{w}z; \theta \mid \lambda) = \int \gamma_1(\Delta z; \theta \mid \lambda, u) f(u; \theta \mid \lambda) du,$$

where $\gamma_1(\Delta z; \theta \mid \lambda, u)$ is Gaussian conditional probability measure of $Z(t)$ given the values λ, u of $A, U, f(u; \theta \mid \lambda)$ the conditional probability density of U (which may contain a linear combination of delta-functions) given the value λ of A , and integration extends over the region of all possible values of U .

From (5.3) we have

$$(5.4) \quad \frac{d\gamma(z; \theta \mid \mu)}{d\gamma(z; \theta \mid \lambda)} = \int \frac{f(v; \theta \mid \mu) dv}{\int \frac{d\gamma_1(z; \theta \mid \lambda, u)}{d\gamma_1(z; \theta \mid \mu, v)} f(u; \theta \mid \lambda) du}$$

and the problem is reduced to calculating Radon-Nikodym derivative of one Gaussian measure with respect to another Gaussian measure.

To calculate this derivative we suppose in addition that for any possible values λ, u of A, U the random function $Z(t)$ is representable by the series

$$(5.5) \quad Z(t) = \sum_{v=1}^{\infty} Z_v x_v(t) \quad (t \in T),$$

where $\{x_v(t)\}$ is some set of functions independent of λ, u (but which may depend on the numerical parameter θ). As such set of functions the set of coordinate functions of a canonical expansion of the random function $Z(t)$ for some specific values λ_0, u_0 of λ, u may be used. Let $\{\Omega^{(v)}\}$ be the set of linear functionals satisfying the biortho-

$$(5.6) \quad \Omega^{(v)} x_\mu = \delta_{v\mu}.$$

Such set may be determined, for instance, using the procedure indicated in [30], § 63. If $x_\nu(t)$ represent the coordinate functions of some canonical expansion, then the set $\{\Omega^{(v)}\}$ is automatically determined while finding this canonical expansion [30], §§ 60–62.

From (5.5) and (5.6) follows

$$(5.7) \quad Z_v = \Omega^{(v)} Z(t).$$

Formulas (5.5) and (5.7) establish complete equivalence between the random function $Z(t)$, $t \in T$, and the set of random variables $\{Z_v\}$ in the sense that any set of values $\{z_v\}$ of $\{Z_v\}$ determine completely the corresponding sample $z(t)$, $t \in T$ of $Z(t)$ and vice versa.

Let us now find the conditional probability density $f_1(\lambda, u | z)$ of A, U given the value $z(t)$, $t \in T$ of $Z(t)$. We have in virtue of the above proved equivalence between $Z(t)$, $t \in T$, and $\{Z_v\}$

$$(5.8) \quad f_1(\lambda, u | z) = \lim_{n \rightarrow \infty} f_1^{(n)}(\lambda, u | z_1, \dots, z_n).$$

On the other hand

$$(5.9) \quad f_1(\lambda, u | z) = \frac{\alpha_1(\lambda) f(u; \theta | \lambda)}{\int_L \alpha_1(\mu) d\mu \int \frac{d\gamma_1(z; \theta | \mu, v)}{d\gamma_1(z; \theta | \lambda, u)} f(v; \theta | \mu) dv}$$

Hence, once $f_1(\lambda, u | z)$ is found, the derivative $d\gamma_1(z; \theta | \mu, v)/d\gamma_1(z; \theta | \lambda, u)$ will be determined immediately. To determine $f_1^{(n)}(\lambda, u | z_1, \dots, z_n)$ in (5.8) we notice that the joint conditional distribution of Z_1, \dots, Z_n is normal for any λ, u with the expected values

$$(5.10) \quad m_\nu(\lambda, u) = \Omega_\nu^{(v)} m_z(t | \lambda, u)$$

and covariance matrix $K_\mu(\lambda, u)$ whose elements are given by

$$(5.11) \quad k_{\nu\mu}(\lambda, u) = \Omega_\nu^{(v)} \Omega_\mu^{(\mu)} K_z(t, \tau | \lambda, u),$$

where $m_z(t | \lambda, u)$, $K_z(t, \tau | \lambda, u)$ are respectively the conditional expected value and conditional covariance of the random function $Z(t)$ given λ, u (they depend also on θ which is omitted for brevity). Therefore

$$(5.12) \quad f_1^{(n)}(\lambda, u | z_1, \dots, z_n) = c_n \alpha_1(\lambda) f(u; \theta | \lambda) \times$$

$$\times \exp \left\{ -\frac{1}{2} \sum_{\nu, \mu=1}^n k_{\nu\mu}^-(\lambda, u) \Omega_\tau^{(\nu)} \Omega_\tau^{(\mu)} [z(t) - m_z(t | \lambda, u)] [z(\tau) - m_z(\tau | \lambda, u)] - \frac{1}{2} \ln A_n(\lambda, u) \right\},$$

where $k_{\nu\mu}^-(\lambda, u)$ represent the elements of the inverse matrix $K_n^{-1}(\lambda, u)$ (necessarily depending on n), c_n the normalizing constant independent of λ, u , and $A_n(\lambda, u)$ the determinant of the $n \times n$ matrix $\Gamma_n(\lambda, u)$ with the elements

$$(5.13) \quad \gamma_{\mu\nu}(\lambda, u) = \frac{k_{\nu\mu}(\lambda, u)}{\sqrt{(D_\nu D_\mu)}},$$

$\{D_\nu\}$ being an arbitrary sequence of positive numbers.

Now introducing the bilinear operator

$$(5.14) \quad Q_n(\lambda, u) = \sum_{\nu, \mu=1}^n [k_{\nu\mu}^-(\lambda, u) - D_\nu^{-1} \delta_{\nu\mu}] \Omega_\tau^{(\nu)} \Omega_\tau^{(\mu)}$$

and the linear operator

$$(5.15) \quad L_n(\lambda, u) = \sum_{\nu=1}^n \left[\sum_{\mu=1}^n k_{\nu\mu}^-(\lambda, u) \Omega_\tau^{(\mu)} m_z(\tau | \lambda, u) \right] \Omega_\tau^{(\nu)},$$

we rewrite (5.12) in the form

$$(5.16) \quad f_1^{(n)}(\lambda, u | z_1, \dots, z_n) = c_n \alpha_1(\lambda) f(u; \theta | \lambda) \exp \left\{ -\frac{1}{2} Q_n(\lambda, u) z(t) z(\tau) + L_n(\lambda, u) z(t) - \frac{1}{2} \beta_n(\lambda, u) \right\},$$

where

$$(5.17) \quad \beta_n(\lambda, u) = L_n(\lambda, u), m_z(t | \lambda, u) + \ln A_n(\lambda, u).$$

Substituting (5.16) into (5.8), taking the limit with $n \rightarrow \infty$, and comparing the result with (5.9), we find

$$(5.18) \quad \frac{d\gamma_1(z; \theta | \mu, v)}{d\gamma_1(z; \theta | \lambda, u)} = \frac{\exp \left\{ -\frac{1}{2} Q(\mu, v) z(t) z(\tau) + L(\mu, v) z(t) - \frac{1}{2} \beta(\mu, v) \right\}}{\exp \left\{ -\frac{1}{2} Q(\lambda, u) z(t) z(\tau) + L(\lambda, u) z(t) - \frac{1}{2} \beta(\lambda, u) \right\}},$$

where

$$(5.19) \quad Q(\lambda, u) = \lim_{n \rightarrow \infty} Q_n(\lambda, u), \quad L(\lambda, u) = \lim_{n \rightarrow \infty} L_n(\lambda, u),$$

and

$$(5.20) \quad \beta(\lambda, u) = \lim_{n \rightarrow \infty} \beta_n(\lambda, u) = L(\lambda, u), m_z(t | \lambda, u) + \ln A(\lambda, u)$$

with $A(\lambda, u) = \lim_{n \rightarrow \infty} A_n(\lambda, u)$.

368 Now we introduce the random function

$$(5.21) \quad X(t) = \sum_{v=1}^{\infty} V_v x_v(t) \quad (t \in T),$$

where V_v are uncorrelated random variables with zero expected values and variances equal to the respective numbers D_v chosen so that the series

$$(5.22) \quad \sum_{v=1}^{\infty} D_v x_v^2(t)$$

be convergent for all $t \in T$. Then the random function $X(t)$ has bounded covariance

$$(5.23) \quad K_x(t, \tau) = \sum_{v=1}^{\infty} D_v x_v(t) x_v(\tau) \quad (t, \tau \in T).$$

Now it is easy to show that the operators Q and L satisfy respectively the linear equations

$$(5.24) \quad Q(\lambda, u)_{,s} K_x(t, \tau | \lambda, u) K_x(s, \sigma) = K_x(\tau, \sigma) - K_x(\tau, \sigma | \lambda, u) \quad (\tau, \sigma \in T),$$

$$(5.25) \quad L(\lambda, u)_{,t} K_x(t, \tau | \lambda, u) = m_x(\tau | \lambda, u) \quad (\tau \in T).$$

Formula (5.18) determines the Radon-Nikodym derivative of two Gaussian measures, if they are absolutely continuous one with respect to other (i.e. equivalent), as is generally the case. But this formula is also valid, if these two measures are orthogonal, giving in such a case Radon-Nikodym derivative in the form of the delta-functional (which is the generalization of usual Dirac delta-function). From (5.18) follow as special cases some results of papers [31, 32].

Substituting (5.18) into (5.4) formulas (5.1) and (5.2) give again (4.6) with

$$(5.26) \quad g(z; \theta | \lambda) = \int f(u; \theta | \lambda) \exp \left\{ -\frac{1}{2} Q(\lambda, u) z(t) z(\tau) + L(\lambda, u) z(t) - \frac{1}{2} \beta(\lambda, u) \right\} du.$$

Thus in the case of systems with continuous input and discrete output the posterior probability density of $A(\theta)$ is determined by (4.6) and (4.7), $g(z; \theta | \lambda)$ being the functional of $z(t)$ given by (5.26) with operators Q and L satisfying the linear equations (5.24) and (5.25), and the function β determined by (5.20).

To solve equations (5.24) and (5.25) the method of canonical expansions of random functions may be applied in the general case. Representing the random function $Z(t)$, $t \in T$, for given λ, u by some canonical expansion, we find $P(\lambda, u, \sigma)_{,t} = = Q(\lambda, u)_{,s} K_x(s, \sigma)$ and $L(\lambda, u)$ as shown in [30], §§ 135, 136, after which using some canonical expansion of the form (5.21) of the random function $X(t)$, we obtain by the same techniques the operator $Q(\lambda, u)$. In various special cases other techniques can be used (see, for example, [30], §§ 128–133, and [33]).

The most difficult for practical calculations in evaluating $g(z; \theta | \lambda)$ is the finding of the infinite determinant $A(\lambda, u)$. To derive a suitable expression for $A(\lambda, u)$ we put

$$(5.27) \quad \bar{K}_x(t, \tau | \lambda, u) = K_x(t, \tau) + R(t, \tau | \lambda, u).$$

Then taking into account (5.6) and (5.23), formula (5.13) yields

$$(5.28) \quad \gamma_{\nu\mu}(\lambda, u) = \delta_{\nu\mu} + \frac{\Omega_t^{(\nu)} \Omega_t^{(\mu)} R(t, \tau | \lambda, u)}{\sqrt{(D_\nu D_\mu)}}.$$

Using the conventional expansion for $A_n(\lambda, u)$ and passing to the limit with $n \rightarrow \infty$, we obtain in virtue of linearity of the functionals $\Omega^{(\nu)}$

$$(5.29) \quad A(\lambda, u) = 1 + \sum_{p=1}^{\infty} \frac{1}{p!} \Theta_{i_1 i_1}^{i_1 j_1} \dots \Theta_{i_p j_p}^{i_p j_p} \begin{vmatrix} R_{i_1 j_1}(t_1, \tau_1 | \lambda, u) & \dots & R_{i_1 j_p}(t_1, \tau_p | \lambda, u) \\ \dots & \dots & \dots \\ R_{i_p j_1}(t_p, \tau_1 | \lambda, u) & \dots & R_{i_p j_p}(t_p, \tau_p | \lambda, u) \end{vmatrix},$$

where $R_{ij}(t, \tau | \lambda, u)$ represent the elements of the matrix $R(t, \tau | \lambda, u)$, and Θ is the bilinear operator

$$(5.30) \quad \Theta = \sum_{\nu=1}^n D_\nu^{-1} \Omega_t^{(\nu)} \Omega_t^{(\nu)},$$

Θ_{ii}^{ij} acting on a matrix whose elements are numbered by i, j and represent functions of t, τ .

It is easy to show that the operator Θ satisfies the linear equation

$$(5.31) \quad \Theta_{i\alpha} K_x(t, \tau) K_x(s, \sigma) = K_x(\tau, \sigma) \quad (\tau, \sigma \in T).$$

Formula (5.29) is especially convenient, if $R(t, \tau | \lambda, u)$ is small as compared with $K_x(t, \tau)$. In this case only a few first members of the series in (5.29) will suffice. The numbers D_ν can be chosen optimally so that $K_x(t, \tau)$ coincide with the expected value of $K_x(t, \tau | \lambda, u)$ averaged with respect to λ, u .

Equations (5.24), (5.25) and (5.31) represent linear integral equations of the first kind, if $Z(t)$ is a scalar random function, and sets of simultaneous linear integral equations of the first kind, if $Z(t)$ is a vector random function. Finally, (5.24), (5.25) and (5.31) represent sets of linear algebraic equations, if T is a discrete set of values of the argument t (case of discrete systems).

The special case where $Z(t)$ represents the sum of two independent random functions, one of which has normal conditional distribution for any values λ, u of A, U , and other is independent of A, U normally distributed random functions, was first studied by L. P. Syssov [34] who used canonical expansions in another way and obtained for this special case another expression for $A(\lambda, u)$.

In the special case where $K_x(t, \tau | \lambda, u)$ is independent of λ, u , it can be taken as $K_x(t, \tau)$. Then (5.24) and (5.29) give $Q(\lambda, u) = 0$, $A(\lambda, u) = 1$, and we obtain the results formerly derived in [35–37] (see also [30], § 144).

In the case of systems with continuous input and continuous output we obtain from (2.10), (2.2)–(2.5), (4.4) and (4.5) the expression (5.1) with

$$(5.32) \quad \frac{d\sigma(\xi | \lambda, \lambda)}{d\beta(\xi)} = \left[\int_L d\mu \int_L \frac{d\gamma(z; \theta | \mu)}{d\gamma(z; \theta | \lambda)} \times \right. \\ \left. \times \prod_{i=1}^N \frac{d\gamma(z_i; \theta_i | \mu_i)}{d\gamma(z_i; \theta_i | \lambda_i)} \frac{d\delta_T(\tilde{w}_i; \theta_i | z_i, \hat{w}_i, \mu_i)}{d\delta_T(\tilde{w}_i; \theta_i | z_i, \hat{w}_i, \lambda_i)} \alpha(\mu, \bar{\mu}) d\bar{\mu} \right]^{-1}.$$

The Radon-Nikodym derivative

$$d\delta_T(\tilde{w}; \theta | z, \hat{w}, \mu) / d\delta_T(\tilde{w}; \theta | z, \hat{w}, \lambda)$$

can be calculated by the same techniques as $d\gamma(z; \theta | \mu) / d\gamma(z; \theta | \lambda)$ in the case where \tilde{W} represents a random function $\tilde{W}(p)$, $p \in P$, depending on A and some other random vector A , and having normal conditional distribution for any $z(t)$, $t \in T$, $\hat{w}(s)$, $s \in S$, and any possible values λ, a of A, A . In the case of teaching by show δ_T is independent of $\hat{w}(s)$, and P coincides with S . In the case of the ideal teacher δ_T in (5.32) coincides with κ .

Thus in the case of systems with continuous input and output formula (4.6) for the posterior probability density $\omega(\lambda | \xi)$ is valid with $g(z; \theta | \lambda)$ determined by (5.26) and

$$(5.33) \quad p(\tilde{w}; \theta | z, \hat{w}, \lambda) = d_T(\tilde{w}; \theta | z, \hat{w}, \lambda) = \\ = \int f_T(a; \theta | z, \hat{w}, \lambda) \exp \left\{ -\frac{1}{2} Q_T(z, \hat{w}, \lambda, a) \tilde{w}(p) \tilde{w}(q) + \right. \\ \left. + L_T(z, \hat{w}, \lambda, a) \tilde{w}(p) - \frac{1}{2} \beta_T(z, \hat{w}, \lambda, a) \right\} da,$$

where $f_T(a; \theta | z, \hat{w}, \lambda)$ is the conditional probability density of the random vector A given $z(t)$, $t \in T$, $\hat{w}(s)$, $s \in S$, λ ; $Q_T(z, \hat{w}, \lambda, a)$ the bilinear operator satisfying the linear equation

$$(5.34) \quad Q_T(z, \hat{w}, \lambda, a)_{pa} K_{\tilde{w}}(p, \xi | z, \hat{w}, \lambda, a) K_{\tilde{w}}(q, \eta | z, \hat{w}) = \\ = K_{\tilde{w}}(\xi, \eta | z, \hat{w}) - K_{\tilde{w}}(\xi, \eta | z, \hat{w}, \lambda, a) \quad (\xi, \eta \in P),$$

$L_T(z, \hat{w}, \lambda, a)$ the linear operator satisfying the equation

$$(5.35) \quad L_T(z, \hat{w}, \lambda, a)_p K_{\tilde{w}}(p, \xi | z, \hat{w}, \lambda, a) = m_{\tilde{w}}(\xi | z, \hat{w}, \lambda, a) \quad (\xi \in P),$$

and

$$(5.36) \quad \beta_T(z, \hat{w}, \lambda, a) = L_T(z, \hat{w}, \lambda, a)_p m_{\tilde{w}}(p | z, \hat{w}, \lambda, a) + \ln A_T(z, \hat{w}, \lambda, a),$$

$\Delta_T(z, \hat{w}, \lambda, a)$ being the infinite determinant with the elements

$$(5.37) \quad \gamma_{\mu}^T(z, \hat{w}, \lambda, a) = \frac{r_p^{(\nu)} r_q^{(\mu)} K_{\hat{w}}(p, q | z, \hat{w}, \lambda, a)}{\sqrt{(A_{\nu} A_{\mu})}}$$

$m_{\hat{w}}(p | z, \hat{w}, \lambda, a)$, $K_{\hat{w}}(p, q | z, \hat{w}, \lambda, a)$ the conditional expected value and covariance of the random function $\hat{W}(p)$, $r^{(\nu)}$ the linear functionals similar to $\Omega^{(\nu)}$, A_{ν} positive numbers similar to D_{ν} , and $K_{\gamma}(p, q | z, \hat{w})$ the conditional covariance of the random function $Y(p)$ constructed in the same way as $X(t)$. The infinite determinant $\Delta_T(z, \hat{w}, \lambda, a)$ can be expressed by the formula similar to (5.29).

To obtain the expression for $\varrho(\xi, \hat{w})$ suitable to practical calculations, it is necessary to impose certain restrictions on the loss function.

At first we suppose that the loss function $l(w, \hat{w} | \lambda)$ is an ordinary function of the values of the functions $w(s)$ and $\hat{w}(s)$ at a finite set of points s_1, \dots, s_m . In this case formula (4.9) is valid, $k(w; \theta | z, \lambda)$ being the joint conditional probability density of $W(s_1), \dots, W(s_m)$.

Secondly we consider the loss function which is a functional of $w(s)$, $\hat{w}(s)$, of the form

$$(5.38) \quad l(w, \hat{w} | \lambda) = \Phi_s \sigma(s, w(s), \hat{w}(s) | \lambda),$$

where Φ is a linear functional in a space of ordinary functions of s , $s \in S$, and σ a function of s , $w(s)$, $\hat{w}(s)$ which may depend also on the derivatives $w'(s)$, $w''(s)$, ..., $w^{(p)}(s)$, $\hat{w}'(s)$, $\hat{w}''(s)$, ..., $\hat{w}^{(p)}(s)$, p being some positive integer. In this case

$$(5.39) \quad \varrho(\xi, \hat{w}) = \int_L \omega(\lambda | \xi) d\lambda \Phi_s \int_B \sigma(s, w, \hat{w}(s) | \lambda) k(w; s, \theta | z, \lambda) dw,$$

where $k(w; s, \theta | z, \lambda)$ is the joint conditional probability density of $W(s)$, $W'(s)$, ..., $W^{(p)}(s)$ for a given value of the argument s .

To estimate the performance of the learned system only $\varepsilon_{\lambda \xi}(\delta)$ can be used in the general case. The quantities $\varepsilon_{\lambda}(\delta)$ and $\varepsilon(\delta)$ can be calculated practically only approximately. And only in some special cases one may succeed in obtaining exact values of $\varepsilon_{\lambda}(\delta)$ and $\varepsilon(\delta)$, as we shall see in the following example.

To obtain $\varepsilon_{\lambda \xi}(\delta)$, we use (4.9) and (4.10), or (5.39) and the corresponding expression for $\varrho_{\lambda}(z, \hat{w})$:

$$(5.40) \quad \varrho_{\lambda}(z, \hat{w}) = \Phi_s \int_B \sigma(s, w, \hat{w}(s) | \lambda) k(w; s, \theta | z, \lambda) dw.$$

In the special case of recognizing system with continuous input (its output is always discrete) all the formulas of the preceding section are valid, if the input $Z(t)$ is a random function of the type considered, with $f_w(z; \theta | \lambda)$ given by the formula

372 similar to (5.26):

$$(5.41) \quad f_w(z; \theta | \lambda) = \int f(u; \theta | \lambda) \exp \left\{ -\frac{1}{2} Q_w(\lambda, u) z(t) z(\tau) + L_w(\lambda, u) z(t) - \frac{1}{2} \beta_w(\lambda, u) \right\} du,$$

operators Q_w and L_w being determined by the equations similar to (5.24) and (5.25), and the function β_w by the formula similar to (5.20).

Example. Let us consider the problem of the example of the preceding section, supposing that the input $Z(t)$ excites the system continuously in the interval of time $0 \leq t \leq T$ and the optimal estimate of the parameter U is required at any instant s of the interval $s_0 \leq s \leq T$. The output of the Bayes optimal system with complete information about λ is in this case determined by

$$(5.42) \quad W_{\lambda}^*(s) = \frac{D_u H(s) + \lambda}{b(s) D_u + 1} \quad (s_0 \leq s \leq T),$$

where

$$(5.43) \quad H(s) = \int_0^s g(s, t) Z(t) dt, \quad b(s) = \int_0^s g(s, t) \varphi(t) dt,$$

the weighting function $g(s, t)$ being determined by the integral equation

$$(5.44) \quad \int_0^s g(s, t) K_{\lambda}(t, \tau) dt = \varphi(\tau) \quad (0 \leq \tau \leq T).$$

To find the Bayes optimal learning system we suppose that the system receives the teacher output $\tilde{W}(s)$ continuously in the time interval $s_0 \leq s \leq T$ at each cycle of learning. We assume also that the teacher gives the Bayes optimal estimate $W_{\lambda}^*(s)$ of $W(s) = U$ with random error $\tilde{Y}(s)$ representing a normally distributed random function with zero expected value and the covariance $K_{\tilde{Y}}(s, s')$. Using (5.20), (5.24)–(5.26) and (5.33)–(5.36) and taking into account that $A(\lambda, u) = A_1(z, w, \lambda, a) = 1$ in this case, we obtain

$$(5.45) \quad g(z; \theta | \lambda) = c_1 \exp \left\{ -\frac{[b(T) D_u + 1] \lambda^2 - [D_u \eta(T) + \lambda]^2}{2 D_u [b(T) D_u + 1]} \right\},$$

$$(5.46) \quad d_1(\tilde{w}; \theta | z, \lambda) = c_2 \exp \{ [\zeta - \Phi(\eta)] \lambda - \frac{1}{2} \kappa \lambda^2 \},$$

where c_1, c_2 are constants independent of λ , $\eta(s)$ is determined by (5.43) with $Z(t)$ replaced by $z(t)$, and

$$(5.47) \quad \zeta = \int_{s_0}^T h(s) \tilde{w}(s) ds, \quad \Phi(\eta) = \int_{s_0}^T \frac{h(s) \eta(s) ds}{b(s) D_u + 1}, \quad \kappa = \int_{s_0}^T \frac{h(s) ds}{b(s) D_u + 1},$$

the weighting function $h(s)$ being determined by the integral equation

$$(5.48) \quad \int_{s_0}^T h(s) K_{\tilde{Y}}(s, \sigma) ds = \frac{1}{b(\sigma) D_u + 1} \quad (s_0 \leq \sigma \leq T).$$

Then assuming as in the example of the preceding section the normal prior distribution for λ with the expected value m and variance D , we obtain the normal posterior distribution $\omega(\lambda | \xi)$, the posterior expected value $\lambda^*(s)$ and variance $\Delta(s)$ of λ being determined by

$$(5.49) \quad \lambda^*(s) = \Delta(s) \left\{ \frac{\eta(s)}{b(s) D_u + 1} + \frac{1}{b(T) D_u + 1} \sum_{i=1}^N \eta_i(T) + \sum_{i=1}^N [\xi_i - \Phi(\eta_i)] + \frac{m}{D} \right\},$$

$$(5.50) \quad \Delta(s) = \frac{D[b(s) D_u + 1][b(T) D_u + 1]}{[b(s) D_u + 1] \{ (1 + N \times D) [b(T) D_u + 1] + N b(T) D \} + b(s) [b(T) D_u + 1]}.$$

The output of the Bayes optimal learning system is determined by

$$(5.51) \quad W^*(s) = \frac{D_u H(s) + \lambda^*(s)}{b(s) D_u + 1}.$$

Finally, we shall estimate the results of learning. Using (4.9) and (4.10) with quadratic loss function $l(w, w | \lambda) = (w - w)^2$, we see that $\varepsilon_{\lambda \xi}(\delta_{\text{opt}})$ is independent of λ , ξ , and therefore

$$(5.52) \quad \varepsilon(\delta_{\text{opt}}) = \varepsilon_{\lambda}(\delta_{\text{opt}}) = \varepsilon_{\lambda \xi}(\delta_{\text{opt}}) = \frac{\Delta(s)}{D_u [b(s) D_u + 1]}.$$

Thus the fortunate circumstance — the independence of $\varrho(\xi, w^*)$ and $\varrho_{\lambda}(z, w_{\lambda}^*)$ of λ and ξ — gives in our case the possibility to find $\varepsilon_{\lambda}(\delta_{\text{opt}})$ and $\varepsilon(\delta_{\text{opt}})$ without tedious calculations. Such a case is, however, rarely encountered in problems of practice.

(Received August 22, 1967.)

REFERENCES

- [1] Pugachev V. S.: Statistical Problems of Pattern Recognition Theory. The Third All-Union Conference on Automatic Control. Odessa, 1965.
- [2] Pugachev V. S.: A Bayes Approach to the Theory of Learning Systems. The Third Congress of IFAC, London, 20—25 June, 1966.
- [3] Pugachev V. S.: Optimal Algorithms of Learning of Automatic Systems in the Case of a Non-Ideal Teacher. Doklady Akademii Nauk SSSR 172 (1967), 5, 1039—1042 (in Russian).
- [4] Fabian V., Špaček A.: Experience in Statistical Decision Problems. Czechoslovak Math. Journ. 6 (1956), 190—194.
- [5] Špaček A.: An Elementary Experience Problem. Transactions of the Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1956). Prague 1957, 253—258.
- [6] Driml M., Špaček A.: Continuous Random Decision Processes Controlled by Experience. Transactions of the First Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1956). Prague 1957, 43—60.

- [7] Winkelbauer K.: Experience in Games of Strategy and in Statistical Decision. Transactions of the First Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1956). Prague 1957, 297–354.
- [8] Hanš O.: Random Fixed Point Theorems. Transactions of the First Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1956). Prague 1957, 105–125.
- [9] Prouza L.: Bemerkungen zur Lineare Prediktion mittels eines lernenden Filters. Transactions of the First Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1956). Prague 1957, 37–41.
- [10] Driml M., Hanš O.: On Experience Theory Problems. Transactions of the Second Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1959). Prague 1960, 93–111.
- [11] Driml M., Hanš O.: Continuous Stochastic Approximations. Transactions of the Second Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1959). Prague 1960, 113–122.
- [12] Driml M., Nedoma J.: Stochastic Approximations for Continuous Random Processes. Transactions of the Second Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1959). Prague 1960, 145–158.
- [13] Hanš O., Špaček A.: Random Fixed Point Approximation by Differentiable Trajectories. Transactions of the Second Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1959). Prague 1960, 203–213.
- [14] Šefl O.: Filters and Predictors which Adapt Their Values to the Unknown Parameters of the Input Process. Transactions of the Second Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1959). Prague 1960, 597–608.
- [15] Fabian V.: A Stochastic Approximation Method for Finding Optimal Conditions in Experimental Work and in Self-Adaptive Systems. *Aplikace matematiky* 6 (1961), 162–183 (in Czech).
- [16] Fabian V.: A Block-scheme of an Automatic Optimizer. *Aplikace matematiky* 7 (1962), 426–440 (in Czech).
- [17] Fabian V.: A New One-dimensional Stochastic Approximation Method for Finding a Local Minimum of a Function. Transactions of the Third Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1962). Prague 1964, 85–105.
- [18] Šefl O.: Some Problems of Automatic Control Processes with Unknown Characteristics. Transactions of the Third Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1962). Prague 1964, 611–620 (in Russian).
- [19] Gardner L. A.: Adaptive Predictors. Transactions of the Third Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1962). Prague 1964, 123–192.
- [20] Pugachev V. S.: Method of Determining the Optimum System with a Non-linear Dependence of the Observed Function of the Parameters of Signal. Automatic and Remote Control. Proceedings of the First Congress of IFAC (1960). Butterworth, London 1961, 2, 702–705.
- [21] Robbins H.: An Empirical Bayes Approach to Statistics. Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, 1 (1955), 157–164.
- [22] Robbins H.: An Empirical Bayes Approach to Testing Statistical Hypothesis. Review of the International Statistical Institute 31 (1963), 2, 195–208.
- [23] Robbins H.: An Empirical Bayes Approach to Problems of Statistical Decision Theory. *Annals of Math. Statistics* 35 (1964), 1, 1–20.
- [24] Johns H. V.: Non-parametric Empirical Bayes Procedures. *Annals of Math. Statistics* 28 (1957), 649–669.
- [25] Johns H. V.: An Empirical Bayes Approach to Non-parametric Two-way Classification.

- [26] Cooper D. V.: Adaptive Pattern Recognition and Signal Detection Using Stochastic Approximation. *IEEE Transactions on Electronic Computers EC-13* (1964), 3, 306—307.
- [27] Tsyppkin Ya. Z.: Adaptation, Learning and Self-learning in Automatic Systems. *Avtomatika i Telemekhanika* 27 (1966), 1, 23—61 (in Russian).
- [28] Shaikin M. E.: A Bayes Approach to Self-adapting Filters Using a Finite Number of Teaching Samples. The Third All-Union Conference on Automatic Control. Odessa, 1965.
- [29] Aizerman M. A., Braverman E. M., Rozonoer L. I.: Theoretical Foundations of Potential Functions Method in the Problem of Teaching Automata to Classify Input Situations. *Avtomatika i Telemekhanika* 25 (1964), 6, 917—936 (in Russian).
- [30] Pugachev V. S.: *Theory of Random Functions and its Application to Control Problems*. Pergamon Press, 1965 (English translation of Russian book: В. С. Пугачев: Теория случайных функций и ее применение к задачам автоматического управления. Физматгиз, 1957, 1960, 1962).
- [31] Shepp L. A.: Radon-Nikodym Derivatives of Gaussian Measures. *Annals of Math. Statistics* 37 (1966), 2, 321—354.
- [32] Rozanov Yu. A.: On the Density of Gaussian Distributions and Wiener - Hopf Integral Equations. *Teoria Veroyatnostei i ee Primeneniya* 11 (1966), 1, 170—179 (in Russian).
- [33] Pugachev V. S.: Methods for Solving Sets of Integral Equations Occuring in Problems of Optimization of Multidimensional Systems. Proceedings of the Sixth All-Union Conference on Probability Theory and Mathematical Statistics (1960). Vilnius 1962, 233—237 (in Russian).
- [34] Sysoev L. P.: Estimates of Parameters of Signals Modulated by Random Processes. *Avtomatika i Telemekhanika* 26 (1965), 10, 1255—1261 (in Russian).
- [35] Pugachev V. S.: An Effective Method for Finding a Bayes Decision. Transactions of the Second Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (1959). Prague 1960, 531—540 (in Russian).
- [36] Pugachev V. S.: A method for Determining Optimum Systems Using General Bayes Criteria. *IRE Transactions on Circuit Theory CT-7* (1960), 4, 491—505.
- [37] Pugachev V. S.: A Method for Determining an Optimal System Using a General Bayes Criterion. *Izvestia Akademii Nauk SSSR, OTN, Energetika i Avtomatika* (1960), 2, 83—97 (in Russian).

Optimální učící se systémy

V. S. PUGAČEV

V článku se probírají optimální Bayesovy učící se systémy a základní vlastnosti jejich algoritmů z hlediska obecné teorie učících se systémů. Diskutují se obecné koncepte učícího se systému, učitele, forem učení se, typy učitelů atd. Uvažuje se obecný případ reálného učitele řídicího procesu s náhodnými chybami. Případ ideálního učitele, který se dříve uvažoval výlučně, je pojímán jako zvláštní případ. Zavádí se míra odlišnosti učícího se systému od optimálního systému s úplnou informací.

V. S. Pugachev, Corresponding Member of the Academy of Sciences of the USSR, Institute of Control Problems, 81 Profsoyuznaya Street, Moscow V-485. USSR.