

# CONTRIBUTIONS TO THE THEORY OF STATISTICAL ESTIMATION AND TESTING HYPOTHESES<sup>1</sup>

BY ABRAHAM WALD

1. **Introduction.** Let us consider a family of systems of  $n$  variates  $X_1(\theta^{(1)}, \dots, \theta^{(k)}), \dots, X_n(\theta^{(1)}, \dots, \theta^{(k)})$  depending on  $k$  parameters  $\theta^{(1)}, \dots, \theta^{(k)}$ . A system of  $k$  values  $\theta^{(1)}, \dots, \theta^{(k)}$  can be represented in the  $k$ -dimensional parameter space by the point  $\theta$  with the co-ordinates  $\theta^{(1)}, \dots, \theta^{(k)}$ . Denote by  $\Omega$  the set of all possible points  $\theta$ . For any point  $\theta$  of  $\Omega$  we shall denote by  $P(E \in w|\theta)$  the probability that the sample point  $E = (x_1, \dots, x_n)$  falls into the region  $w$  of the  $n$ -dimensional sample space, where  $x_j$  denotes the observed value of the variate  $X_j(\theta)$  ( $j = 1, \dots, n$ ). The distribution  $P(E \in w|\theta)$  is supposed to be known for any point  $\theta$  of  $\Omega$ . In the theory of testing hypotheses and of statistical estimation we have to deal with problems of the following type: A sample point  $E = (x_1, \dots, x_n)$  of the  $n$ -dimensional sample space is given. We know that  $x_j$  is the observed value of  $X_j(\theta)$  but we do not know the parameter point  $\theta$ , and we have to draw inferences about  $\theta$  by means of the sample point observed. The assumption that  $\theta$  belongs to a certain subset  $\omega$  of  $\Omega$  is called a hypothesis. We shall deal in this paper with the following general problem: Let us consider a system  $S$  of subsets of  $\Omega$ . Denote by  $H_\omega$  the hypothesis corresponding to the element  $\omega$  of  $S$ , and by  $H_S$  the system of all hypotheses corresponding to all elements of  $S$ . We have to decide by means of the observed sample point  $E$  which hypothesis of the system  $H_S$  should be accepted. That is to say for each  $H_\omega$  we have to determine a region of acceptance  $M_\omega$  in the  $n$ -dimensional sample space. The hypothesis  $H_\omega$  will be accepted if and only if the sample point  $E$  falls in the region  $M_\omega$ .  $M_\omega$  and  $M_{\omega'}$  are disjoint if  $\omega \neq \omega'$ . The statistical problem is the question as to how the system  $M_S$  of all regions  $M_\omega$  should be chosen.

The problem in this formulation is very general. It contains the problems of testing hypotheses and of statistical estimation treated in the literature.<sup>2</sup> For instance if we want to test the hypothesis  $H_\omega$  corresponding to a certain subset  $\omega$  of  $\Omega$ , the system of hypotheses  $H_S$  consists only of the two hypotheses  $H_\omega$  and  $H_{\bar{\omega}}$  where  $\bar{\omega}$  denotes the subset of  $\Omega$  complementary to  $\omega$ . If we want to estimate  $\theta$  by a unique point, then  $S$  is the system of all points of  $\Omega$ . In the theory of confidence intervals we estimate one of the parameter co-ordinates  $\theta^{(1)}, \dots, \theta^{(k)}$ ,

<sup>1</sup> Research under a grant-in-aid from the Carnegie Corporation of New York.

<sup>2</sup> See, for instance, J. Neyman, "Outline of a Theory of Statistical Estimation Based on the Classical Theory of Probability," *Phil. Transactions of the Royal Society*, London, Vol. 231 (1937), pp. 333-380.

say  $\theta^{(1)}$ , by an interval. In this case  $S$  is a certain system of subsets  $\omega$  of the following type:  $\omega$  is the set of all points  $\theta = (\theta^{(1)}, \dots, \theta^{(k)})$  for which  $\theta^{(1)}$  lies in a certain interval  $[a, b]$ . The problem in our formulation covers also cases which, as far as I know, have not yet been treated. Consider for instance 3 subsets  $\omega_1, \omega_2$  and  $\omega_3$  of  $\Omega$  such that the sum of them is equal to  $\Omega$ . It may be that we are interested only to know in which of the subsets  $\omega_1, \omega_2, \omega_3$  the unknown parameter point lies. In this case the system of hypotheses  $H_s$  consists only of the 3 hypotheses  $H_{\omega_1}, H_{\omega_2}$  and  $H_{\omega_3}$ . Cases like this might be of practical interest.

For the determination of the "best" system (in a certain sense) of regions of acceptance we shall use methods and principles which are closely related to those of the Neyman-Pearson theory of testing hypotheses. In the Neyman-Pearson theory two types of error are considered. Let  $\theta = \theta_1$  be the hypothesis to be tested, where  $\theta_1$  denotes a certain point of the parameter space. Denote this hypothesis by  $H_1$  and the hypothesis  $\theta \neq \theta_1$  by  $\bar{H}$ . The type I error is that which is made by rejecting  $H_1$  when it is true. The type II error is made by accepting  $H_1$  when it is false. The fundamental principle in the Neyman-Pearson theory can be formulated as follows: among all critical regions (regions of rejection of  $H_1$ , i.e. regions of acceptance of  $\bar{H}$ ) for which the probability of type I error is equal to a given constant  $\alpha$ , we have to choose that region for which the probability of type II error is a minimum. The difficulty which arises here lies in the circumstance that the probability of type II error depends on the true parameter point  $\theta$ . That is to say, if the critical region is given the probability of type II error will be a function of the true parameter point  $\theta$ . Since we do not know the true parameter point  $\theta$ , we want to have a critical region which minimizes the probability of type II error with respect to any possible alternative hypothesis  $\theta = \theta_2 \neq \theta_1$ . If such a common best critical region exists, then the problem is solved. But such cases are rather exceptional. If a common best critical region does not exist, Neyman and Pearson consider unbiased critical regions of different types,<sup>3</sup> which minimize the type II error locally, that is to say with respect to alternative hypotheses in the neighborhood of the hypothesis considered. In this paper we develop methods for the determination of a system of regions of acceptance taking in account type II errors also relative to alternative hypotheses not lying in the neighborhood of the hypothesis to be tested.

**2. Some Definitions.** Let us denote by  $\Omega$  the set of all possible parameter points  $\theta$  and by  $S$  a system of subsets of  $\Omega$ . If  $\rho$  denotes the sum of the elements of a subset  $\sigma$  of  $S$ , then we shall denote  $\Sigma M_\omega$  by  $M_\rho$ , where  $M_\omega$  denotes the

<sup>3</sup> J. Neyman and E. S. Pearson: *Statistical Research Memoirs*, Volumes I and II. The authors consider also unbiased regions of type  $A_1$  for which the probability of type II error with respect to every alternative hypothesis is not greater than for any other unbiased region of the same size. However regions of type  $A_1$  do not always exist (the existence of such regions has been proved for a special but important class of cases).

region of acceptance of  $H_\omega$  and the summation is to be taken over all elements  $\omega$  of  $\sigma$ .

*Definition 1.* Denote by  $M_S$  and  $M'_S$  two different systems of regions of acceptance corresponding to the same system  $H_S$  of hypotheses. The systems  $M_S$  and  $M'_S$  are said to be equivalent if for each point  $\theta$  of  $\Omega$  and for every  $\rho$  which is a sum of elements of  $S$  which does not contain  $\theta$ , the equation

$$P(E \in M'_\rho | \theta) = P(E \in M_\rho | \theta)$$

holds, where  $M'_\rho$  denotes the region according to the system  $M'_S$  and  $M_\rho$  denotes the region according to the system  $M_S$ .

*Definition 2.* Denote by  $M_S$  and  $M'_S$  two different systems of regions of acceptance corresponding to the same system of hypotheses. The system  $M'_S$  is said to be absolutely better than the system  $M_S$  if they are not equivalent and if for each  $\theta$  and for every  $\rho$  which is a sum of elements of  $S$  which does not contain  $\theta$  the inequality

$$P(E \in M'_\rho | \theta) \leq P(E \in M_\rho | \theta)$$

holds.

*Definition 3.* A system  $M_S$  of regions of acceptance is said to be admissible if no absolutely better system of regions exists.

**3. The problem of the choice of  $M_S$ .** The choice of  $M_S$  will in general be affected by the following two circumstances:

(1) We do not attribute the same importance to each error. For instance the acceptance of the hypothesis that  $\theta$  lies in a certain interval  $I$  has in general more serious consequences if  $\theta$  is far from  $I$  than if  $\theta$  is near to  $I$ . The choice of  $M_S$  will in general depend on the relative importance of the different possible errors.

(2) In some cases we have a priori more confidence that the true parameter point lies in a certain interval  $I$  than in some other cases. The choice of  $M_S$  will in general be affected also by this fact. Let us illustrate this by an example. We have two coins, a new and an old one and we want to test for both coins whether the probability  $p$  of tossing head is equal to  $\frac{1}{2}$ . Let us assume that we make 100 tosses with each of the coins and we get head 40 times in each case. Since we have a priori no very great confidence that the old coin is unbiased, the fact that head occurred only 40 times will suffice to reject the hypothesis that for the old coin  $p = \frac{1}{2}$ . But in the case of the new coin, having much greater a priori confidence that it is unbiased, we shall perhaps not reject the hypothesis  $p = \frac{1}{2}$  and we shall rather assume that a somewhat improbable event occurred. That is to say, we do not choose the same critical region in both cases due to the fact that our a priori confidence for  $p = \frac{1}{2}$  is in the case of the new coin greater than in the case of the old one.

In order to study the dependence of the choice of  $M_S$  on the two circumstances

mentioned, let us introduce a weight function for the possible errors and an a priori probability distribution for the unknown parameter  $\theta$ . The weight function  $W(\theta, \omega)$  is a real valued non-negative function defined for all points  $\theta$  of  $\Omega$  and for all elements  $\omega$  of  $S$ , which expresses the relative importance of the error committed by accepting  $H_\omega$  when  $\theta$  is true. If  $\theta$  is contained in  $\omega$ ,  $W(\theta, \omega)$  is, of course, equal to zero. The question as to how the form of the weight function  $W(\theta, \omega)$  should be determined, is not a mathematical or statistical one. The statistician who wants to test certain hypotheses must first determine the relative importance of all possible errors which will entirely depend on the special purposes of his investigation. If that is done, we shall in general be able to give a more satisfactory answer to the question as to how the system of regions of acceptance should be chosen. In many cases, especially in statistical questions concerning industrial production, we are able to express the importance of an error in terms of money, that is to say, we can express the loss caused by the error considered in terms of money. We shall also say that  $W(\theta, \omega)$  is the loss caused by accepting  $H_\omega$  when  $\theta$  is true.

The situation regarding the introduction of an a priori probability distribution of  $\theta$  is entirely different. First, the objection can be made against it, as Neyman has pointed out, that  $\theta$  is merely an unknown constant and not a variate, hence it makes no sense to speak of the probability distribution of  $\theta$ . Second, even if we may assume that  $\theta$  is a variate, we have in general no possibility of determining the distribution of  $\theta$  and any assumptions regarding this distribution are of hypothetical character. On account of these facts the determination of the system of regions of acceptance should be independent of any a priori probability considerations. The "best" system of regions of acceptance, which we shall define later, will depend only on the weight function of the errors. The reason why we introduce here a hypothetical probability distribution of  $\theta$  is simply that it proves to be useful in deducing certain theorems and in the calculation of the best system of regions of acceptance.

Let us denote by  $f(\theta)$  a distribution function of  $\theta$ . For the sake of simplicity let us assume that the probability density of the distribution  $P(E \in \omega | \theta)$  exists in any point  $E$  of the sample space for any  $\theta$  and denote it by  $p(E | \theta)$ . The expected value of the loss is given by

$$(1) \quad I = \int_M \int_\Omega W(\theta, \omega_E) p(E | \theta) df(\theta) dE$$

where  $\omega_E$  denotes the element of  $S$  corresponding to  $E$  (that is to say,  $\omega_E$  is that element of  $S$  for which  $E$  is a point of the region of acceptance  $M_{\omega_E}$ ), and the integral is to be taken over the product of the sample space  $M$  with the parameter space  $\Omega$ . The expected value  $I$  of the loss depends on the system  $M_S$  of regions of acceptance. The system  $M_S$  for which  $I$  becomes a minimum, can be regarded as the best system of regions relative to the given weight function and to the given a priori distribution of  $\theta$ .

One can easily show the following: If  $M'_S$  is an absolutely better system of regions (in sense of the definition 2) than the system  $M_S$ , then for any weight

function  $w(\theta, \omega)$  and for any a priori distribution  $f(\theta)$  the expected value  $I'$  of the loss corresponding to  $M'_s$  is less than the expected value  $I$  of the loss corresponding to  $M_s$ . (For some exceptional weight and a priori distribution functions  $I'$  may be equal to  $I$ .)

Hence we can give the following rule: *We have to choose an admissible system of regions of acceptance.*

Now let us consider the question whether besides admissibility further restrictions upon the choice of  $M_s$  can be made. In order to see this, let us consider two admissible systems of regions  $M_s$  and  $M'_s$  which are not equivalent. One can easily show that there exist two weight functions  $W_1(\theta, \omega)$ ,  $W_2(\theta, \omega)$  and two a priori distributions  $f_1(\theta)$  and  $f_2(\theta)$  such that for  $W_1(\theta, \omega)$  and  $f_1(\theta)$  the expected value of the loss corresponding to  $M_s$  is less than that corresponding to  $M'_s$ , and for  $W_2(\theta, \omega)$  and  $f_2(\theta)$  the expected value of the loss corresponding to  $M_s$  is greater than that corresponding to  $M'_s$ . Hence no absolute criteria can be given as to which of the systems  $M_s$  and  $M'_s$  should be chosen. In order to be able to make further restrictions upon the choice of  $M_s$ , we have to make assumptions regarding the form of the weight function. We shall deal with this question in section 6.

**4. Calculation of admissible systems of regions.** As we have seen, we have to choose an admissible system of regions. The question arises as to how we can find admissible systems of regions.

Provided that  $p(E | \theta)$  is continuous in  $E$  and  $\theta$  jointly, one can easily show that  $M'_s$  is an admissible system of regions if there exists a bounded, uniformly continuous and everywhere positive (except if  $\theta$  is contained in  $\omega$ ) weight function  $W(\theta, \omega)$  and an a priori distribution  $f(\theta)$  such that every open subset of  $\Omega$  has a positive probability and the expected value of the loss

$$(2) \quad I(M_s) = \int_M \int_{\Omega} W(\theta, \omega_E) p(E | \theta) df(\theta) dE.$$

becomes a minimum for  $M_s = M'_s$ . ( $\omega_E$  denotes that element of  $S$  for which  $M_{\omega_E}$  contains  $E$ ). In fact if there existed an absolutely better system  $M''_s$  of regions, then  $I(M''_s)$  would be less than  $I(M'_s)$  in contradiction to our assumption that  $I(M_s)$  becomes a minimum for  $M_s = M'_s$ .

In order to obtain an admissible system  $M_s$  we may choose any bounded, uniformly continuous and everywhere positive (except if  $\theta$  is contained in  $\omega$ ) weight function  $W(\theta, \omega)$  and any arbitrary a priori distribution  $f(\theta)$  (subject only to the condition that every open subset of  $\Omega$  should have a positive probability) and then the system  $M_s$  which makes

$$I(M_s) = \int_M \int_{\Omega} W(\theta, \omega_E) p(E | \theta) df(\theta) dE$$

a minimum is an admissible one. In order to determine  $M_s$  we have only to determine for each  $E$  the corresponding element  $\omega_E$  of  $S$ . Let us consider the integral

$$I_E = \int_{\Omega} W(\theta, \omega) p(E | \theta) df(\theta).$$

The integral  $I_E$  is for a fixed  $E$  only a function of  $\omega$ . It is obvious that  $\omega_E$  must be that element of  $S$  for which  $I_E$  becomes a minimum.

**5. Admissible systems  $M_S$  and the Neyman-Pearson best critical regions.**

Let us consider the case that the system  $H_S$  of hypotheses consists only of the following two hypotheses: 1)  $\theta = \theta_0$  where  $\theta_0$  is a certain point of  $\Omega$ . 2)  $\theta$  belongs to the set complementary to  $\theta_0$ . Let us denote by  $\omega_1$  the set consisting only of the point  $\theta_0$ , and by  $\omega_2$  the set complementary to  $\omega_1$ .  $S$  consists in this case only of two elements  $\omega_1$  and  $\omega_2$ . The system  $M_S$  of regions consists of two regions of acceptance  $M_{\omega_1}$  and  $M_{\omega_2}$  corresponding to the hypotheses  $H_{\omega_1}$  and  $H_{\omega_2}$ . If a common best critical region in the sense of Neyman-Pearson exists and if  $M_S$  is admissible, then  $M_{\omega_2}$  is obviously a common best critical region. This leads to the following remarkable conclusion: If a common best critical region exists and if the system  $M_S$  of regions consisting of the two regions  $M_{\omega_1}$  and  $M_{\omega_2}$  minimizes the expectation of the loss (formula 2) for a weight function and for an a priori distribution subject to some weak conditions mentioned in paragraph 4, then  $M_{\omega_2}$  is a common best critical region. That is to say, the form of the weight function and of the a priori distribution affects only the size of the region  $M_{\omega_2}$  but it will always be a common best critical region.

**6. The choice of  $M_S$  if a weight function is given.** We shall now consider the case in which a weight function  $W(\theta, \omega)$  is given and we shall deal with the question as to how  $M_S$  in this case is to be chosen.

If the parameter point is an unknown constant and if  $\theta$  denotes the true parameter point, then the expected value of the loss is given by

$$(3) \quad r(\theta) = \int_M W(\theta, \omega_E) p(E | \theta) dE$$

where the integration is to be taken over the whole sample space  $M$  and  $H_{\omega_E}$  denotes the hypothesis accepted if  $E$  is the observed sample point. That is to say  $\omega_E$  is that element of  $S$  for which  $E$  is contained in the region of acceptance  $M_{\omega_E}$ . We shall call the expression (3) the risk of accepting a false hypothesis if  $\theta$  is the true parameter point. Since we do not know the true parameter point  $\theta$ , we shall have to study the risk  $r(\theta)$  as a function of  $\theta$ . We shall call this function the risk function. The form of the risk function depends on the system  $M_S$  of regions and on the form of the weight function. In order to express this fact, we shall denote the risk function corresponding to the system  $M_S$  and to the weight function  $W(\theta, \omega)$  also by

$$r[\theta | M_S, W(\theta, \omega)].$$

*Definition 4.* Denote by  $M_S$  and  $M'_S$  two systems of regions of acceptance corresponding to the same system  $H_S$  of hypotheses. We shall say that  $M_S$  and  $M'_S$  are equivalent relative to the weight function  $W(\theta, \omega)$  if the risk function

$r[\theta | M_s, W(\theta, \omega)]$  is identically equal to the risk function  $r[\theta | M'_s, W(\theta, \omega)]$ , that is to say if for each point  $\theta$ ,

$$r[\theta | M'_s, W(\theta, \omega)] = r[\theta | M_s, W(\theta, \omega)].$$

*Definition 5.* Denote by  $M_s$  and  $M'_s$  two systems of regions corresponding to the same system  $H_s$  of hypotheses. We shall say that  $M_s$  is uniformly better than  $M'_s$  relative to the weight function  $W(\theta, \omega)$  if  $M_s$  and  $M'_s$  are not equivalent and for each  $\theta$

$$r[\theta | M_s, W(\theta, \omega)] \leq r[\theta | M'_s, W(\theta, \omega)].$$

*Definition 6.* A system  $M_s$  of regions of acceptance is said to be admissible relative to the weight function  $W(\theta, \omega)$  if no uniformly better system of regions exists relative to the weight function considered.

It is obvious that we have to choose a system  $M_s$  of regions which is admissible relative to the weight function considered.

There exist in general many systems  $M_s$  which are admissible relative to the weight function given. The question arises as to how can we distinguish among them. Denote by  $r_{M_s}$  the maximum of the risk function corresponding to the system  $M_s$  of regions and to the given weight function. If we do not take into consideration a priori probabilities of  $\theta$ , then it seems reasonable to choose that system  $M_s$  for which  $r_{M_s}$  becomes a minimum. We shall see in section 8 that the system  $M_s$  for which  $r_{M_s}$  becomes a minimum has some important properties which justify the distinction of this particular system of regions among all admissible systems.

*Definition 7.* We shall call an admissible system  $M'_s$  of regions for which  $r_{M'_s}$  becomes a minimum a best system of regions of acceptance relative to the weight function given.<sup>4</sup>

Now we shall have to deal with the question of determining a best system  $M_s$  of regions and what special properties this system  $M_s$  has.

**7. Reduction of the problem to the case when the system  $H_s$  of hypotheses is the system of all simple hypotheses.** A hypothesis  $H_\omega$  is said to be a simple hypothesis if  $\omega$  contains exactly one point of the parameter space  $\Omega$ . We assume that each element  $\omega$  of  $S$  is a closed subset of  $\Omega$ . Hence the power of  $S$  is not greater than the power of the continuum and therefore we can always set up a correspondence between the elements  $\omega$  of  $S$  and the points  $\theta$  of  $\Omega$  such that to each point  $\theta$  corresponds a certain element  $\omega_\theta$  of  $S$  and to each element  $\omega$  of  $S$  at least one point  $\theta$  exists for which  $\omega_\theta = \omega$ . For instance if  $S$  consists of the two elements  $\omega_1$  and  $\omega_2$  then we can set up a correspondence as follows: the element  $\omega_\theta$  of  $S$  corresponding to  $\theta$  is  $\omega_1$  if  $\theta$  is contained in  $\omega_1$  and  $\omega_2$  otherwise.

---

<sup>4</sup> As we shall see later (Theorem 3), the best system of regions is uniquely determined if some regularity conditions are fulfilled.

If  $\Omega$  is one dimensional and  $S$  is the system of all intervals of a certain length  $\epsilon$  then we can define the interval  $\omega_\theta$  corresponding to  $\theta$  as the interval of which the initial point is  $\theta$  and the terminal point  $\theta + \epsilon$ .

Let us denote the weight function by  $W(\theta, \omega)$  defined for all values of  $\theta$  and for all elements  $\omega$  of  $S$ . Consider the system  $H_{\bar{s}}$  of all simple hypotheses and the following weight function

$$(4) \quad W(\theta, \bar{\theta}) = W(\theta, \omega_{\bar{\theta}})$$

where  $\theta$  denotes the true parameter point and  $\bar{\theta}$  denotes the estimated point. A system  $M_{\bar{s}}$  of regions of acceptance for  $H_{\bar{s}}$  is given by a vector function  $\bar{\theta}(E)$  of the observations such that to each point  $E = (x_1, \dots, x_n)$  of the sample space  $M$  corresponds a certain point  $\bar{\theta}(E)$  of the parameter space. For each point  $\theta_0$  the region  $M_{\theta_0}$  of the acceptance of the hypothesis  $\theta = \theta_0$  is given by the equation  $\bar{\theta}(E) = \theta_0$ . We shall call the function  $\bar{\theta}(E)$  an estimate of  $\theta$ , the system of regions  $M_{\bar{s}}$  is uniquely determined by the estimate. We shall call  $\bar{\theta}(E)$  a best estimate relative to a given weight function if the system of regions determined by  $\bar{\theta}(E)$  is a best system of regions relative to the weight function considered.

Let us denote by  $\bar{\theta}(E)$  a best estimate of  $\theta$  relative to the weight function  $W(\theta, \bar{\theta})$  defined in (4). A best system  $M_s$  of regions of acceptance in the original problem can obviously be obtained in the following way: Denote by  $\omega$  an element of  $S$ . The region  $M_\omega$  of acceptance of the hypothesis  $H_\omega$  consists of the points  $E$  for which

$$\omega_{\bar{\theta}(E)} = \omega.$$

Hence we can restrict our considerations to the case when the system of hypotheses is the system of all simple hypotheses. We shall deal with the problem of how a best estimate of  $\theta$  can be found and what properties this estimate has.

**8. Some theorems concerning the best estimate.** In order to study the properties of a best estimate  $\bar{\theta}(E)$  it is useful to consider hypothetical a priori distributions of  $\theta$ . We shall especially consider point distributions of  $\theta$ , that is to say, distributions where a finite number of points  $\theta_1, \dots, \theta_s$  of the parameter space  $\Omega$  exist such that the probability of any subset of  $\Omega$  not containing any of the points  $\theta_1, \dots, \theta_s$  is zero. If  $\theta_1, \dots, \theta_s$  are given, a point distribution is characterized by a vector  $\rho = (\rho_1, \dots, \rho_s)$  where  $\rho_i$  denotes the probability of  $\theta_i$  and  $\sum \rho_i = 1$ .

If  $\theta(E)$  denotes an estimate of  $\theta$  and if  $f(\theta)$  denotes a distribution function of  $\theta$  then the expected value of the loss, that is to say the expected value of the weight function  $W[\theta, \theta(E)]$  is obviously given by

$$(5) \quad \int_M \int_\Omega W[\theta, \theta(E)] p(E | \theta) df(\theta) dE$$

where  $p(E | \theta)$  denotes the probability density in  $E$  if  $\theta$  is the true parameter point and the integration is to be taken over the product of the sample space  $M$  and parameter space  $\Omega$ .



Let us assume that for every sample point  $E$  there exists a parameter point  $\theta_f(E)$  such that the expression

$$(6) \quad \int_{\Omega} W(\theta, \bar{\theta})p(E|\theta) df(\theta)$$

becomes a minimum with respect to  $\bar{\theta}$  for  $\bar{\theta} = \theta_f(E)$ . We shall call the estimate  $\theta_f(E)$  a minimum risk estimate with respect to the distribution  $f(\theta)$ , since also the expression (5) becomes a minimum for the estimate  $\theta_f(E)$ .

We shall make the following assumptions:

*Assumption 1.* The parameter space is a bounded and closed subset of the  $k$ -dimensional Euclidean space.

*Assumption 2.* The weight function  $W(\theta, \bar{\theta})$  is continuous in  $\theta$  and  $\bar{\theta}$  jointly.

*Assumption 3.* The probability density  $p(E|\theta)$  is continuous in  $E$  and  $\theta$  jointly. That is to say if  $\lim E_i = E$  and  $\lim \theta_i = \theta$  then  $\lim p(E_i|\theta_i) = p(E|\theta)$ .

*Assumption 4.* For any distribution  $f(\theta)$  of  $\theta$  there exists at most one minimum risk estimate  $\theta_f(E)$ .<sup>5</sup>

*Assumption 5.* If  $f(\theta)$  and  $f'(\theta)$  denote two different point distributions of  $\theta$  and if  $\theta_f(E)$  and  $\theta_{f'}(E)$  are minimum risk estimates corresponding to  $f(\theta)$  and  $f'(\theta)$  respectively, then  $\theta_f(E)$  is not identically equal to  $\theta_{f'}(E)$ .

The assumptions 1-5, with addition of an assumption 6 which we shall formulate later, enables us to deduce important properties of the best estimate  $\hat{\theta}(E)$ . First we shall prove some Lemmas by means of the assumptions 1-5.

**LEMMA 1.** *For any a priori distribution  $f(\theta)$  of  $\theta$  there exists exactly one minimum risk estimate  $\theta_f(E)$ .*

According to Assumption 2  $W(\theta, \bar{\theta})$  is continuous. Since the parameter space  $\Omega$  is compact on account of Assumption 1,  $W(\theta, \bar{\theta})$  is uniformly continuous. According to Assumption 3  $p(E|\theta)$  is continuous; hence for any fixed sample point  $E$ ,  $p(E|\theta)$  is bounded. From these facts it follows easily that the expression (6) is a continuous function of  $\bar{\theta}$  for any fixed sample point  $E$ . Hence there exists at least one parameter point  $\theta_f(E)$  such that (6) becomes a minimum for  $\bar{\theta} = \theta_f(E)$ . Since, according to Assumption 4, at most one parameter point exists for which (6) becomes a minimum, Lemma 1 is proved.

If a distribution  $f(\theta)$  of  $\theta$  is given then the distribution of each of the components  $\theta^{(1)}, \dots, \theta^{(k)}$  of  $\theta$  can be found. Denote by  $Q_j$  the set of real numbers which are discontinuities of the distribution of the component  $\theta^{(j)}$  ( $j = 1, \dots, k$ ) and form the set  $Q = Q_1 + \dots + Q_k$ . As is well known,  $Q$  is at most denumerable. A  $k$ -dimensional interval  $J$  of the parameter space given by

$$a_j \leq \theta^{(j)} \leq b_j \quad (j = 1, \dots, k)$$

is called a continuity interval of the distribution  $f(\theta)$  if no  $a_j$  and no  $b_j$  belongs to  $Q$ . A sequence  $\{f_n(\theta)\}$  of distributions is said to be convergent towards the

---

<sup>5</sup> As will be shown in Section 10, Assumption 4 is not as restrictive as it would appear. It will be satisfied in the great majority of practical cases.

distribution  $f(\theta)$ , i.e. in symbols  $\lim f_n(\theta) = f(\theta)$ , if for any continuity interval  $J$  of  $f(\theta)$  the probability of  $J$  corresponding to the distribution  $f_n(\theta)$  converges with increasing  $n$  towards the probability of  $J$  corresponding to the distribution  $f(\theta)$ .

LEMMA 2. *If  $\{f_n(\theta)\}$  ( $n = 1, \dots, \text{ad inf.}$ ) denotes a sequence of distributions, then there exists a subsequence  $\{f_{n_m}(\theta)\}$  ( $m = 1, \dots, \text{ad inf.}$ ) which converges towards a distribution.*

As is well known, there exists a completely additive set function  $P(\omega)$  defined for all Borel measurable subsets  $\omega$  of  $\Omega$  and a subsequence  $\{n_m\}$  of  $\{n\}$ , such that for any continuity interval  $J$  of  $P(\omega)$  the probability of  $J$  corresponding to the distribution  $f_{n_m}(\theta)$  converges with increasing  $m$  towards  $P(J)$ . Since  $\Omega$  is bounded, there exists a continuity interval  $J$  such that for all  $n$  the probability of  $J$  according to  $f_n(\theta)$  is equal to 1. Hence  $P(\Omega) = 1$ , that is to say,  $P(\omega)$  is a probability set function which proves Lemma 2.

LEMMA 3. *If  $\{f_n(\theta)\}$  ( $n = 1, \dots, \text{ad inf.}$ ) denotes a sequence of distributions which converges towards the distribution  $f(\theta)$  and if  $\lim E_n = E$  then*

$$\lim_{n \rightarrow \infty} \theta_{f_n}(E_n) = \theta_f(E),$$

where  $\theta_{f_n}(E)$  denotes the minimum risk estimate corresponding to  $f_n(\theta)$  and  $\theta_f(E)$  denotes the minimum risk estimate corresponding to  $f(\theta)$ .

If  $\{\varphi_n(\theta)\}$  denotes a sequence of real valued functions which converges uniformly towards a continuous function  $\varphi(\theta)$  then

$$(7) \quad \lim \int_{\Omega} \varphi_n(\theta) df_n(\theta) = \int_{\Omega} \varphi(\theta) df(\theta).$$

Since  $\{\varphi_n(\theta)\}$  converges uniformly towards  $\varphi(\theta)$ , (7) is obviously true if

$$\lim \int_{\Omega} \varphi(\theta) df_n(\theta) = \int_{\Omega} \varphi(\theta) df(\theta)$$

holds. The latter equality follows easily from the fact that  $\Omega$  is compact.

Consider a subsequence  $\{n_m\}$  of  $\{n\}$  such that  $\lim_{m \rightarrow \infty} \theta_{f_{n_m}}(E_{n_m})$  exists. Denote this limit by  $\theta^*$ . In order to prove Lemma 3, we have only to show that  $\theta^* = \theta_f(E)$ . If  $\theta_f(E) \neq \theta^*$  then on account of Assumption 4

$$(8) \quad \int_{\Omega} W[\theta, \theta_f(E)] p(E | \theta) df(\theta) < \int_{\Omega} W(\theta, \theta^*) p(E | \theta) df.$$

$W(\theta, \bar{\theta})$  is uniformly continuous since  $\Omega$  is compact. On account of Assumption 3 also  $p(E | \theta)$  is uniformly continuous in the product of  $\Omega$  with a bounded subset of the sample space. Hence

$$W[\theta, \theta_{f_{n(m)}}(E_{n_m})] p(E_{n_m} | \theta)$$

converges uniformly in  $\theta$  towards

$$W(\theta, \theta^*) p(E | \theta)$$

and we have on account of (7) and (8)

$$(9) \quad \lim_{m \rightarrow \infty} \int_{\Omega} W[\theta, \theta_{f_n(m)}(E_{n_m})]p(E_{n_m} | \theta) df_{n_m} = \int_{\Omega} W(\theta, \theta^*)p(E | \theta) df > \int_{\Omega} W[\theta, \theta_f(E)]p(E | \theta) df,$$

and

$$(10) \quad \lim_{m \rightarrow \infty} \int_{\Omega} W[\theta, \theta_f(E)]p(E | \theta) df_{n_m} = \int_{\Omega} W[\theta, \theta_f(E)]p(E | \theta) df.$$

From (9) and (10) it follows that there exists a positive  $\delta$  such that for sufficiently large  $m$

$$\int_{\Omega} W[\theta, \theta_{f_n(m)}(E_{n_m})]p(E_{n_m} | \theta) df_{n_m} > \int_{\Omega} W[\theta, \theta_f(E)]p(E | \theta) df_{n_m} + \delta.$$

Since the sequence of functions  $\{p(E_n | \theta)\}$  converges uniformly in  $\theta$  towards  $p(E | \theta)$ , we have for sufficiently large  $m$

$$\int_{\Omega} W[\theta, \theta_{f_n(m)}(E_{n_m})]p(E_{n_m} | \theta) df_{n_m} > \int_{\Omega} W[\theta, \theta_f(E)]p(E_{n_m} | \theta) df_{n_m}.$$

But this is a contradiction, since  $\theta_{f_n}(E)$  is a minimum risk estimate. Hence the assumption  $\theta^* \neq \theta_f(E)$  is proved to be an absurdity. This proves Lemma 3.

LEMMA 4. *To each positive  $\epsilon$  a bounded and closed subset  $M_{\epsilon}$  of the sample space  $M$  can be given such that*

$$\int_{M_{\epsilon}} p(E | \theta) dE \geq 1 - \epsilon$$

for every point  $\theta$  of the parameter space  $\Omega$ .

Let us assume that Lemma 4 is not true and we shall deduce a contradiction. Denote by  $M_{\nu}$  ( $\nu = 1, 2, \dots$ , ad inf.) the sphere in the sample space  $M$  whose center is the origin and whose radius is equal to  $\nu$ . Since Lemma 4 is supposed to be not true, to each  $\nu$  there exists a parameter point  $\theta_{\nu}$  such that

$$(11) \quad \int_{M_{\nu}} p(E | \theta_{\nu}) dE < 1 - \epsilon \quad (\nu = 1, \dots, \text{ad inf.}).$$

Since  $\Omega$  is compact, there exists a subsequence  $\{\theta_{\nu_{\mu}}\}$  of the sequence  $\{\theta_{\nu}\}$  such that  $\lim_{\mu \rightarrow \infty} \theta_{\nu_{\mu}}$  exists. Denote  $\lim \theta_{\nu_{\mu}}$  by  $\theta$ . Since

$$\int_M p(E | \theta) dE = 1$$

there exists a positive integer  $\nu'$  such that

$$\int_{M_{\nu'}} p(E | \theta) dE > 1 - \frac{\epsilon}{2}.$$

On account of Assumption 3 we get easily

$$\lim_{\mu \rightarrow \infty} \int_{M, \mu} p(E | \theta_{\nu, \mu}) dE = \int_{M, \nu} p(E | \theta) dE.$$

Hence for sufficiently large  $\mu$  we get

$$\int_{M, \nu, \mu} p(E | \theta_{\nu, \mu}) dE \geq \int_{M, \nu} p(E | \theta_{\nu, \mu}) dE > 1 - \epsilon,$$

in contradiction to (11). This proves Lemma 4.

For any estimate  $\theta(E)$  we shall call the integral

$$r(\theta) = \int_M W[\theta, \theta(E)] p(E | \theta) dE$$

the risk function of the estimate  $\theta(E)$ . The value of the risk function  $r(\theta)$  is for any  $\theta$  equal to the expected value of the loss (of the weight function) if  $\theta$  is the true parameter point.

LEMMA 5. To any positive  $\eta$  a positive  $\delta$  can be given such that for any estimate  $\theta(E)$  and for any pair  $\theta, \theta'$  of parameter points whose Euclidean distance is less than  $\delta$  the inequality

$$|r(\theta) - r(\theta')| = \left| \int_M W[\theta, \theta(E)] p(E | \theta) dE - \int_M W[\theta', \theta(E)] p(E | \theta') dE \right| < \eta$$

holds.

Since  $W(\theta, \bar{\theta})$  is uniformly continuous, to any  $\epsilon > 0$  a positive  $\delta$  can be given such that for any pair of points  $\theta, \theta'$  whose Euclidean distance is less than  $\delta$  the relation

$$(12) \quad |W(\theta, \bar{\theta}) - W(\theta', \bar{\theta})| < \epsilon$$

holds for every  $\bar{\theta}$ . On account of Assumption 3  $\delta$  can be chosen in such a way that also the inequality

$$(13) \quad |p(E | \theta) - p(E | \theta')| < \epsilon$$

is satisfied for any sample point  $E$  of a bounded subset  $M'$  of  $M$  and for any pair  $\theta, \theta'$  whose Euclidean distance is less than  $\delta$ .

Since  $W(\theta, \bar{\theta})$  is continuous and  $\Omega$  is compact,  $W(\theta, \bar{\theta})$  must be bounded. Denote by  $A$  an upper bound of  $W(\theta, \bar{\theta})$ . According to Lemma 4 there exists a bounded and closed subset  $M'$  of the sample space  $M$  such that

$$\int_{M'} p(E | \theta) dE \geq 1 - \frac{\eta}{2A} \text{ for any } \theta.$$

It is obvious that

$$\left| \int_{M-M'} W[\theta, \theta(E)] p(E | \theta) dE - \int_{M-M'} W[\theta', \theta(E)] p(E | \theta') dE \right| \leq \frac{\eta}{2}.$$

In order to prove Lemma 5 we have only to show that

$$(14) \quad \left| \int_{M'} W[\theta, \theta(E)]p(E | \theta) dE - \int_{M'} W[\theta', \theta(E)]p(E | \theta') dE \right| < \frac{\eta}{2}.$$

On account of (12) and (13), (14) is certainly true for sufficiently small  $\epsilon$ . Hence Lemma 5 is proved.

LEMMA 6. *If the sequence  $\{f_n(\theta)\}$  of distributions converges towards the distribution  $f(\theta)$  and if  $r_{f_n}(\theta)$  denotes the risk function of the minimum risk estimate  $\theta_{f_n}(E)$  then  $\{r_{f_n}(\theta)\}$  converges uniformly towards the risk function  $r_f(\theta)$  of the minimum risk estimate  $\theta_f(E)$ .*

According to Lemma 4 to any positive  $\epsilon$  a bounded and closed subset  $M_\epsilon$  of  $M$  can be given such that

$$(15) \quad \int_{M_\epsilon} p(E | \theta) dE \geq 1 - \epsilon$$

for every  $\theta$ . From Lemma 3 it follows easily that  $\{\theta_{f_n}(E)\}$  converges uniformly towards  $\theta_f(E)$  in  $M_\epsilon$ . Hence

$$\lim_{n \rightarrow \infty} \int_{M_\epsilon} W[\theta, \theta_{f_n}(E)]p(E | \theta) dE = \int_{M_\epsilon} W[\theta, \theta_f(E)]p(E | \theta) dE$$

holds for every  $\theta$  and for every positive  $\epsilon$ . Since  $W(\theta, \bar{\theta})$  is bounded and  $\epsilon$  can be chosen arbitrarily small, we get on account of (15) that

$$\lim_{n \rightarrow \infty} \int_M W[\theta, \theta_{f_n}(E)]p(E | \theta) dE = \int_M W[\theta, \theta_f(E)]p(E | \theta) dE,$$

that is to say

$$\lim r_{f_n}(\theta) = r_f(\theta).$$

The uniformity of the convergence follows easily from Lemma 5.

In the following argument we shall consider an arbitrary but fixed system of  $s$  parameter points  $\theta_1, \dots, \theta_s$ , and point distributions such that no point  $\theta \neq \theta_1, \dots, \theta_s$  has positive probability. Such a point distribution is characterized by a vector  $\rho = (\rho_1, \dots, \rho_s)$  where  $\rho_i$  denotes the probability of  $\theta_i$  ( $i = 1, \dots, s$ ) and  $\sum \rho_i = 1$ . The points  $\theta_1, \dots, \theta_s$  are kept constant and only  $\rho$  will vary. Hence if we speak about different distributions  $\rho = (\rho_1, \dots, \rho_s)$ ,  $\rho' = (\rho'_1, \dots, \rho'_s)$  they are always related to the same points  $\theta_1, \dots, \theta_s$  unless we state explicitly the contrary.

LEMMA 7. *If  $\rho = (\rho_1, \dots, \rho_s)$  and  $\rho' = (\rho_1 + \Delta\rho_1, \dots, \rho_s + \Delta\rho_s)$  denote two different distributions then*

$$\sum_{i=1}^s [(\lambda - 1)\rho_i + \lambda\Delta\rho_i][r_i(\rho') - r_i(\rho)] < 0$$

holds for any positive  $\lambda$ , where

$$r_i(\rho) = \int_M W[\theta_i, \theta_\rho(E)]p(E | \theta_i) dE \quad (i = 1, \dots, s),$$

$$r_i(\rho') = \int_M W[\theta_i, \theta_{\rho'}(E)]p(E | \theta_i) dE,$$

and  $\theta_\rho(E)$  and  $\theta_{\rho'}(E)$  denote the minimum risk estimates corresponding to  $\rho$  and  $\rho'$  respectively.

We have

$$\sum_i (\rho_i + \Delta\rho_i)r_i(\rho) = \int_M \sum_i W[\theta_i, \theta_\rho(E)]\rho'_i p(E | \theta_i) dE = I_1$$

and

$$\sum_i (\rho_i + \Delta\rho_i)r_i(\rho') = \int_M \sum_i W[\theta_i, \theta_{\rho'}(E)]\rho'_i p(E | \theta_i) dE = I_2.$$

Since  $\theta_{\rho'}(E)$  is the minimum risk estimate corresponding to  $\rho'$ , we have  $I_1 \geq I_2$ . We shall show that  $I_1 > I_2$ . According to Assumption 5  $\theta_\rho(E)$  is not identically equal to  $\theta_{\rho'}(E)$ . Hence there exists a point  $E'$  such that  $\theta_\rho(E') \neq \theta_{\rho'}(E')$ . On account of Assumption 4

$$\Sigma W[\theta_i, \theta_\rho(E')] \rho'_i p(E' | \theta_i) > \Sigma W[\theta_i, \theta_{\rho'}(E')] \rho'_i p(E' | \theta_i).$$

From Lemma 3 it follows that  $\theta_\rho(E)$  and  $\theta_{\rho'}(E)$  are continuous functions of  $E$ . Hence there exists a positive  $\delta$  and a sphere  $s$  with center in  $E'$  such that

$$\Sigma W[\theta_i, \theta_\rho(E)] \rho'_i p(E | \theta_i) > \Sigma W[\theta_i, \theta_{\rho'}(E)] \rho'_i p(E | \theta_i) + \delta$$

for every point  $E$  of  $S$ . Since  $\theta_{\rho'}(E)$  is the minimum risk estimate corresponding to  $\rho'$  we have

$$\Sigma W[\theta_i, \theta_\rho(E)] \rho'_i p(E | \theta_i) \geq \Sigma W[\theta_i, \theta_{\rho'}(E)] \rho'_i p(E | \theta_i)$$

for every point  $E$  outside  $S$ . Hence  $I_1 > I_2$  that is to say

$$(16) \quad \Sigma(\rho_i + \Delta\rho_i)r_i(\rho) > \Sigma(\rho_i + \Delta\rho_i)r_i(\rho').$$

Analogously we get

$$(17) \quad \Sigma\rho_i r_i(\rho) < \Sigma\rho_i r_i(\rho').$$

Multiplying (16) by an arbitrary positive value  $\lambda$  and subtracting (17) we get

$$\Sigma[\lambda(\rho_i + \Delta\rho_i) - \rho_i]r_i(\rho) > \Sigma[\lambda(\rho_i + \Delta\rho_i) - \rho_i]r_i(\rho').$$

Hence

$$\Sigma[(\lambda - 1)\rho_i + \lambda\Delta\rho_i][r_i(\rho') - r_i(\rho)] < 0.$$

Let us denote for any  $\rho$  the maximum of the numbers

$$r_1(\rho), \dots, r_s(\rho)$$

by  $r(\rho)$ . We shall call a distribution  $\rho$  for which  $r(\rho)$  becomes a minimum, a risk-minimizing distribution. We shall say that the risk-minimizing distribution  $\rho = (\rho_1, \dots, \rho_s)$  is not degenerate if  $\rho_1 > 0, \dots, \rho_s > 0$ . Otherwise we shall say that  $\rho$  is degenerate.

LEMMA 8. *There exists at least one risk-minimizing distribution  $\rho$ .*

From Lemma 6 it follows that  $r_1(\rho), \dots, r_s(\rho)$  are continuous functions of  $\rho$ . Hence also  $r(\rho)$  is continuous. Since the set of all possible distributions  $\rho$  is bounded and closed, there must be at least one distribution  $\rho$  for which  $r(\rho)$  becomes a minimum.

LEMMA 9. *If  $\rho = (\rho_1, \dots, \rho_s)$  denotes a risk-minimizing distribution which is not degenerate then*

$$r_1(\rho) = r_2(\rho) = \dots = r_s(\rho).$$

Let us assume that there are two integers  $i$  and  $j$ , for instance 1 and 2, such that  $r_1(\rho) < r_2(\rho)$ . We shall deduce a contradiction from this assumption. Let us consider two different distributions  $\rho' = (\rho'_1, \dots, \rho'_s)$  and  $\rho'' = (\rho''_1, \dots, \rho''_s)$  where  $\rho''_1 > 0$ . Hence at least one of the quantities

$$(\rho'_1 - \rho''_1), \dots, (\rho'_s - \rho''_s)$$

is unequal to zero. Since  $\sum \rho'_i = \sum \rho''_i = 1$ , also at least one of the quantities

$$(\rho'_2 - \rho''_2), \dots, (\rho'_s - \rho''_s)$$

must be unequal to zero. On account of Lemma 7 we have

$$\sum_{i=1}^s [(\lambda - 1)\rho'_i + \lambda(\rho''_i - \rho'_i)][r_i(\rho'') - r_i(\rho')] < 0.$$

If we put  $\lambda = \frac{\rho'_1}{\rho''_1}$  we get

$$\sum_{i=2}^s \left[ \left( \frac{\rho'_1}{\rho''_1} - 1 \right) \rho'_i + \frac{\rho'_1}{\rho''_1} (\rho''_i - \rho'_i) \right] [r_i(\rho'') - r_i(\rho')] < 0.$$

Hence at least one of the quantities

$$r_2(\rho'') - r_2(\rho'), \dots, r_s(\rho'') - r_s(\rho')$$

must be unequal to zero.

Since  $\rho_1 > 0$ , there exists a closed sphere  $S_\rho$  with center at  $\rho$  such that for any point  $\rho'$  of  $S_\rho$   $\rho'_1 > 0$ . Hence for any two different points  $\rho'$  and  $\rho''$  of  $S_\rho$  at least one of the quantities

$$r_2(\rho'') - r_2(\rho'), \dots, r_s(\rho'') - r_s(\rho')$$

is unequal to zero. Denote by  $\bar{S}_\rho$  the projection of  $S_\rho$  on the  $s - 1$  dimensional space given by  $\rho_1 = 0$ . Consider the transformation according to which the image of the point  $\bar{\rho}' = (\rho'_2, \dots, \rho'_s)$  of  $\bar{S}_\rho$  is the point  $\bar{q}(\bar{\rho}') = [r_2(\rho'), \dots, r_s(\rho')]$ . It is obvious that the images of two different points of  $\bar{S}_\rho$  are different.

Since  $r_i(\rho)$  ( $i = 1, \dots, s$ ) is continuous, the transformation is continuous and therefore topological. Denote the image of  $\bar{S}_\rho$  by  $\bar{R}_\rho$ . Since  $\bar{\rho} = (\rho_2, \dots, \rho_s)$  is an interior point of  $\bar{S}_\rho$ , according to the Brouwer-Jordan theorem<sup>6</sup> on domain invariance the image  $\bar{q}(\bar{\rho}) = [r_2(\rho), \dots, r_s(\rho)]$  of  $\bar{\rho}$  must also be an interior point of  $\bar{R}_\rho$ . Hence for sufficiently small  $\epsilon > 0$  the point

$$t(\epsilon) = [r_2(\rho) - \epsilon, \dots, r_s(\rho) - \epsilon]$$

is contained in  $\bar{R}_\rho$ . Denote by  $\bar{\rho}(\epsilon) = [\rho_2(\epsilon), \dots, \rho_s(\epsilon)]$  the point of  $\bar{S}_\rho$  whose image is  $t(\epsilon)$ . It is obvious that

$$(18) \quad \lim_{\epsilon \rightarrow 0} \bar{\rho}(\epsilon) = \bar{\rho} = (\rho_2, \dots, \rho_s).$$

Consider the point  $\rho(\epsilon)$  of  $S_\rho$  whose projection is  $\bar{\rho}(\epsilon)$  that is to say  $\rho(\epsilon)$  has the co-ordinates  $1 - \Sigma \bar{\rho}_i(\epsilon), \bar{\rho}_2(\epsilon), \dots, \bar{\rho}_s(\epsilon)$ . From (18) it follows that also

$$(19) \quad \lim_{\epsilon \rightarrow 0} \rho(\epsilon) = \rho = (\rho_1, \rho_2, \dots, \rho_s).$$

Since  $r_1[\rho(\epsilon)], \dots, r_s[\rho(\epsilon)]$  are continuous functions of  $\epsilon$  and since  $r_1(\rho) < r_2(\rho)$ , for sufficiently small  $\epsilon$  the maximum of the numbers

$$r_1[\rho(\epsilon)], r_2[\rho(\epsilon)] = r_2(\rho) - \epsilon, \dots, r_s[\rho(\epsilon)] = r_s(\rho) - \epsilon$$

is certainly smaller than the maximum  $r(\rho)$  of the numbers

$$r_1(\rho), \dots, r_s(\rho),$$

in contradiction to our assumption that  $\rho$  is a risk minimizing distribution. Hence the assumption  $r_1(\rho) < r_2(\rho)$  is proved to be an absurdity and Lemma 9 is proved.

In the previous arguments we have considered an arbitrary but fixed system of  $s$  parameter points  $\theta_1, \dots, \theta_s$  and all distributions  $\rho$  were related to these points. In the following arguments we shall vary the points  $\theta_1, \dots, \theta_s$  and therefore we shall have to state the parameter points to which the distribution  $\rho$  is related.

Let us consider a sequence  $\{\theta_\nu\}$  ( $\nu = 1, \dots, \text{ad inf.}$ ) of parameter points which is dense in  $\Omega$ . We say that a subset  $\omega$  of  $\Omega$  is dense in  $\Omega$  if for each point  $\theta$  of  $\Omega$  any arbitrarily small open neighborhood of  $\theta$  contains at least one point of  $\omega$ . Since  $\Omega$  is compact, a sequence  $\{\theta_\nu\}$  which is dense in  $\Omega$  certainly exists. Let us consider the first  $s$  points  $\theta_1, \dots, \theta_s$  of the sequence  $\{\theta_\nu\}$ . According to Lemma 8 there exists for any  $s$  a risk-minimizing distribution  $\rho(s) = [\rho_1(s), \dots, \rho_s(s)]$  related to  $\theta_1, \dots, \theta_s$ .

*Assumption 6.* There exists a sequence  $\{\theta_s\}$  ( $s = 1, \dots, \text{ad inf.}$ ) of parameter points which is dense in  $\Omega$  and such that for almost any  $s$ <sup>7</sup> the risk-minimizing

<sup>6</sup> See for instance Alexandroff and Hopf, *Topologie*, Berlin 1935, p. 396.

<sup>7</sup> By "almost any  $s$ " we understand "for all  $s$  greater than a sufficiently large integer."



distribution  $\rho(s) = [\rho_1(s), \dots, \rho_s(s)]$  related to the first  $s$  points  $\theta_1, \dots, \theta_s$ , is not degenerate.

LEMMA 10. Denote by  $\{\theta_s\}$  ( $s = 1, 2, \dots$ , ad inf.) a sequence of parameter points for which the conditions of Assumption 6 are fulfilled. Denote by  $\rho(s) = [\rho_1(s), \dots, \rho_s(s)]$  the risk-minimizing distribution related to the first  $s$  points  $\theta_1, \dots, \theta_s$ . Then there exists a non-negative constant  $c$  such that for any arbitrarily small positive  $\epsilon$  the inequality

$$c - \epsilon \leq \int_M W[\theta, \theta_{\rho(s)}(E)]p(E | \theta) dE \leq c + \epsilon$$

holds identically in  $\theta$  for almost every  $s$ . That is to say the risk function of the minimum risk estimate  $\theta_{\rho(s)}(E)$  lies entirely between  $c - \epsilon$  and  $c + \epsilon$  for almost every  $s$ .

Denote the risk function

$$\int_M W[\theta, \theta_{\rho(s)}(E)]p(E | \theta) dE$$

of the estimate  $\theta_{\rho(s)}(E)$  by  $r(\theta, s)$ . First we shall prove that there exists a sequence  $\{c_s\}$  ( $s = 1, \dots$ , ad inf.) of non-negative numbers such that for every  $\epsilon > 0$  the inequality

$$(20) \quad c_s - \epsilon \leq r(\theta, s) \leq c_s + \epsilon$$

holds for almost every  $s$ . In fact to any positive  $\eta$  a positive integer  $s_\eta$  can be given such that for any  $s > s_\eta$  the points  $\theta_1, \dots, \theta_s$  are  $\eta$ -dense in  $\Omega$ . That is to say every point  $\theta$  of  $\Omega$  lies in a sphere with radius  $\eta$  and center in one of the points  $\theta_1, \dots, \theta_s$ . Since for sufficiently large  $s$   $\rho(s)$  is not degenerate, we have on account of Lemma 9 for sufficiently large  $s$

$$(21) \quad r(\theta_1, s) = \dots = r(\theta_s, s) = c_s.$$

Since for sufficiently large  $s$   $\theta_1, \dots, \theta_s$  is  $\eta$ -dense in  $\Omega$ , we get easily from Lemma 5 that (20) holds for any positive  $\epsilon$  for almost every  $s$ .

In order to prove Lemma 10 we have only to show that  $\lim_{s \rightarrow \infty} c_s$  exists and is finite. First we see that for no estimate  $\theta(E)$  can the corresponding risk function

$$r(\theta) = \int_M W[\theta, \theta(E)]p(E | \theta) dE$$

lie entirely below  $r(\theta, s)$  that is to say

$$(22) \quad r(\theta) < r(\theta, s)$$

cannot hold for any  $\theta$ . In fact if (22) were true for a certain estimate  $\theta(E)$  then

$$\begin{aligned} \sum \rho_i(s)r(\theta_i) &= \int_M \sum W[\theta_i, \theta(E)]\rho_i(s)p(E | \theta_i) dE < \sum \rho_i(s)r(\theta_i, s) \\ &= \int_M \sum W[\theta_i, \theta_{\rho(s)}(E)]\rho_i(s)p(E | \theta_i) dE, \end{aligned}$$

which is not possible since  $\theta_{\rho(s)}(E)$  is a minimum risk estimate. Hence (22) cannot hold for any  $\theta$ . From this fact follows easily that  $\lim c_s$  exists and is finite. This proves Lemma 10.

LEMMA 11. Denote  $f(\theta)$  a distribution of  $\theta$  and let  $\theta_f(E)$  be the corresponding minimum risk estimate. If  $\theta(E)$  denotes an arbitrary estimate then

$$r(\theta) \equiv r_f(\theta)$$

if  $\theta_f(E) \neq \theta(E)$  only in a set of measure 0, and

$$\int_{\Omega} r(\theta) df(\theta) > \int_{\Omega} r_f(\theta) df(\theta)$$

if  $\theta_f(E) \neq \theta(E)$  in a set of positive measure.  $r(\theta)$  denotes the risk function of  $\theta(E)$  and  $r_f(\theta)$  denotes the risk function of  $\theta_f(E)$ .

If  $\theta_f(E) \neq \theta(E)$  only in a set of measure zero, then we have obviously  $r(\theta) \equiv r_f(\theta)$ . Consider the case that  $\theta_f(E) \neq \theta(E)$  in a set  $M'$  of positive measure. According to Assumption 4 we have

$$\int_{\Omega} W[\theta, \theta(E)] p(E | \theta) df(\theta) > \int_{\Omega} W[\theta, \theta_f(E)] p(E | \theta) df(\theta)$$

for any point  $E$  of  $M'$ . Since

$$\int_{\Omega} W[\theta, \theta(E)] p(E | \theta) df(\theta) = \int_{\Omega} W[\theta, \theta_f(E)] p(E | \theta) df(\theta)$$

for any other point  $E$  of the sample space  $M$ , we get

$$\begin{aligned} \int_{\Omega} r(\theta) df &= \int_M \int_{\Omega} W[\theta, \theta(E)] p(E | \theta) df dE \\ &> \int_M \int_{\Omega} W[\theta, \theta_f(E)] p(E | \theta) df dE = \int_{\Omega} r_f(\theta) df. \end{aligned}$$

Hence Lemma 11 is proved.

We are now able to prove some theorems about the best estimate  $\bar{\theta}(E)$  relative to a given weight function. An estimate  $\bar{\theta}(E)$  is a best estimate according to our definition 7, if the maximum of the risk function of  $\bar{\theta}(E)$  is less than or equal to the maximum of the risk function of any other estimate  $\theta(E)$  and if  $\bar{\theta}(E)$  is an admissible estimate (that is to say there exists no estimate  $\theta(E)$  such that the risk function  $r(\theta)$  of  $\theta(E)$  is not identical to the risk function  $\bar{r}(\theta)$  of  $\bar{\theta}(E)$  and in every point  $\theta$   $\bar{r}(\theta) \geq r(\theta)$ ).

THEOREM 1. If  $\bar{\theta}(E)$  is a best estimate and if the Assumptions 1-6 are fulfilled then the risk function  $\bar{r}(\theta)$  of  $\bar{\theta}(E)$  is constant, that is to say

$$\bar{r}(\theta) \equiv c.$$

According to Assumption 6 there exists a sequence  $\{\theta_s\}$  ( $s = 1, \dots, \text{ad inf.}$ ) of parameter points such that  $\{\theta_s\}$  is dense in  $\Omega$  and for almost every  $s$  the risk-

minimizing distribution  $\rho(s)$  related to  $\theta_1, \dots, \theta_s$  is not degenerate. On account of Lemma 10 there exists a non-negative constant  $c$  such that for any  $\epsilon > 0$  the inequality

$$(23) \quad c - \epsilon \leq r(\theta, s) \leq c + \epsilon$$

holds for almost every  $s$ .  $r(\theta, s)$  denotes the risk function of the estimate  $\theta_{\rho(s)}(E)$ . According to Lemma 2 there exists a subsequence  $\{s_n\}$  ( $n = 1, \dots$ , ad inf.) of integers such that the sequence  $\{\rho(s_n)\}$  of distributions converges towards a distribution  $f(\theta)$ . From Lemma 6 it follows that

$$\lim_{n \rightarrow \infty} r(\theta, s_n) = r_f(\theta)$$

where  $r_f(\theta)$  denotes the risk function of the minimum risk estimate  $\theta_f(E)$ . On account of (23) we have

$$r_f(\theta) \equiv c.$$

From Lemma 11 it follows that for any other estimate  $\theta(E)$  either

$$r(\theta) \equiv r_f(\theta) \equiv c$$

or

$$\int_{\Omega} r(\theta) df > \int_{\Omega} r_f(\theta) df,$$

where  $r(\theta)$  denotes the risk function of  $\theta(E)$ . In the latter case there exists at least one point  $\theta$  for which  $r(\theta) > r_f(\theta)$ . Hence  $\theta_f(E)$  is a best estimate. If  $\bar{\theta}(E)$  is also a best estimate, we get on account of Lemma 11 that  $\bar{\theta}(E)$  can differ from  $\theta_f(E)$  only in a set of measure 0 and the risk function of  $\bar{\theta}(E)$  is identically equal to  $c$ . Hence we have proved Theorem 1 and also the following Theorems 2-3:

**THEOREM 2.** *If the Assumptions 1-6 are fulfilled there exists a distribution  $f(\theta)$  of  $\theta$  such that the corresponding minimum risk estimate  $\theta_f(E)$  is a best estimate.*

**THEOREM 3.** *If Assumptions 1-6 are fulfilled and  $\bar{\theta}(E)$ ,  $\theta^*(E)$  are best estimates, then  $\bar{\theta}(E) = \theta^*(E)$  almost everywhere and the corresponding risk functions are identically equal.*

Now we shall prove (without making the Assumptions 1-6)

**THEOREM 4.** *If  $W(\theta, \bar{\theta})$  and  $p(E | \theta)$  are continuous and  $\Omega$  is compact, and if  $f(\theta)$  denotes a distribution of  $\theta$  such that any open set has a positive probability, then the minimum risk estimate  $\theta_f(E)$  is a best estimate if its risk function  $r_f(\theta)$  is identically equal to a constant.*

Let  $r_f(\theta)$  be identically equal to  $c$  and consider an arbitrary estimate  $\theta(E)$ . Since  $W(\theta, \bar{\theta})$  and  $p(E | \theta)$  are continuous and  $\Omega$  is compact, the risk function  $r(\theta)$  of  $\theta(E)$  is a continuous function of  $\theta$ . Since  $\theta_f(E)$  is a minimum risk estimate we have

$$(24) \quad \int_{\Omega} r(\theta) df \geq \int_{\Omega} r_f(\theta) df = c.$$

In order to prove Theorem 4, we have to show that either

$$(25) \quad r(\theta) \equiv c$$

or there exists a point  $\theta'$  such that

$$(26) \quad r(\theta') > c.$$

If (25) does not hold there exists a point  $\theta^*$  such that  $r(\theta^*) \neq c$ . If  $r(\theta^*) > c$  our statement is proved. Consider the case  $r(\theta^*) < c$ . On account of the continuity of  $r(\theta)$  there exists a positive  $\delta$  and an open neighborhood  $U$  of  $\theta^*$  such that

$$r(\theta) < c - \delta$$

for every  $\theta$  in  $U$ . Since  $\int_U df$  is assumed to be positive, the inequality (24) can hold only if there exists at least a point  $\theta'$  for which  $r(\theta') > c$ . This proves Theorem 4.

**9. Determination of the best estimate  $\bar{\theta}(E)$  for a certain class of distributions  $p(E | \theta)$ .** In this paragraph we shall prove two theorems which enable us to calculate very easily the best estimate  $\bar{\theta}(E)$  for a certain special but important class of distributions.

The risk function of an estimate  $\bar{\theta}(E)$  is given by

$$r(\theta) = \int_M W[\theta, \bar{\theta}(E)] p(E | \theta) dE,$$

where the integration is to be taken over the whole sample space  $M$ . We consider the integral equation

$$(27) \quad \int_M W[\theta, \bar{\theta}(E)] p(E | \theta) dE \equiv c,$$

where  $c$  denotes an arbitrary constant. If we can find an estimate  $\bar{\theta}(E)$  which satisfies (27) for a certain  $c$  and which is an admissible estimate relative to the weight function considered, then  $\bar{\theta}(E)$  is certainly a best estimate. If Assumptions 1-6 are fulfilled, an admissible estimate satisfying (27) certainly exists. As we shall see, a best estimate can very easily be determined by the above procedure if the conditions in the following theorem 5 are fulfilled.

**THEOREM 5.** *Let us assume that the following conditions are fulfilled:*

I. *The parameter space  $\Omega$  is one dimensional and  $\theta$  can take any real value from  $-\infty$  to  $+\infty$ .*

II. *The probability density  $p(E | \theta)$  depends only on the differences  $x_1 - \theta, \dots, x_n - \theta$ , that is to say  $p(E | \theta) = p(x_1 - \theta, \dots, x_n - \theta)$ , where  $x_1, \dots, x_n$  denote the co-ordinates of  $E$ .*

III. *The value of the weight function depends only on the difference  $u = \theta - \bar{\theta}$  and is uniformly continuous in  $u$ .*

IV. For any value  $\bar{\theta}$  and for any sample point  $E$  the integral

$$(28) \quad \psi(\bar{\theta}, E) = \int_{-\infty}^{+\infty} W(\theta - \bar{\theta})p(E | \theta) d\theta$$

has a finite value.

V. For every  $E$  there exists a finite value  $\theta'(E)$  such that  $\psi(\bar{\theta}, E)$  becomes a minimum for  $\bar{\theta} = \theta'(E)$ .

Then there exists an estimate  $\bar{\theta}(E)$  such that for any  $E$ ,  $\psi(\bar{\theta}, E)$  becomes a minimum for  $\bar{\theta} = \bar{\theta}(E)$  and  $\bar{\theta}(E'') - \bar{\theta}(E') = \lambda$  for any  $E' = (x'_1, \dots, x'_n)$  and  $E'' = (x''_1, \dots, x''_n)$  for which  $x''_1 - x'_1 = \dots = x''_n - x'_n = \lambda$ . An estimate with these properties is a best estimate.

Let us consider two sample points  $E' = (x'_1, \dots, x'_n)$  and  $E'' = (x''_1, \dots, x''_n)$  such that  $x''_1 - x'_1 = \dots = x''_n - x'_n = \lambda$ . From the conditions II and III follows that if  $\psi(\bar{\theta}, E')$  becomes a minimum for  $\bar{\theta} = \theta_1$ , then  $\psi(\bar{\theta}, E'')$  becomes a minimum for  $\bar{\theta} = \theta_2 = \theta_1 + \lambda$ . Hence there exists an estimate  $\bar{\theta}(E) = \bar{\theta}(x_1, \dots, x_n)$  such that for any  $E$ ,  $\psi(\bar{\theta}, E)$  becomes a minimum for  $\bar{\theta} = \bar{\theta}(E)$  and  $\bar{\theta}(E'') - \bar{\theta}(E') = \lambda$  if  $x''_1 - x'_1 = \dots = x''_n - x'_n = \lambda$ . We shall show that such an estimate  $\bar{\theta}(E)$  is a best estimate. First we shall show that the risk function

$$r(\theta) = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} W[\theta - \bar{\theta}(E)] p(x_1 - \theta, \dots, x_n - \theta') dx_1 \dots dx_n$$

is constant. Let us consider two arbitrary parameter values  $\theta'$  and  $\theta''$ . Then we have

$$r(\theta') = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} W[\theta' - \bar{\theta}(E)] p(x_1 - \theta', \dots, x_n - \theta') dx_1 \dots dx_n,$$

$$r(\theta'') = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} W[\theta'' - \bar{\theta}(E)] p(x_1 - \theta'', \dots, x_n - \theta'') dx_1 \dots dx_n.$$

Making in the second integral the transformation

$$y_1 = x_1 - (\theta'' - \theta'), \dots, y_n = x_n - (\theta'' - \theta'),$$

we get

$$\begin{aligned} r(\theta'') &= \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} W\{\theta'' - \bar{\theta}[y_1 + (\theta'' - \theta'), \dots, y_n \\ &\quad + (\theta'' - \theta')]\} p(y_1 - \theta', \dots, y_n - \theta') dy_1 \dots dy_n \\ &= \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} W[\theta' - \bar{\theta}(y_1, \dots, y_n)] p(y_1 - \theta', \dots, y_n - \theta') dy_1 \dots dy_n. \end{aligned}$$

Hence  $r(\theta') = r(\theta'')$  and our statement that  $r(\theta)$  is constant is proved. In order to prove Theorem 5, we have only to show that  $\bar{\theta}(E)$  is an admissible estimate. For this purpose let us consider an arbitrary estimate  $\theta^*(E)$  and

denote the corresponding risk function by  $r^*(\theta)$ . Since  $\bar{\theta}(E)$  minimizes the integral (28), we have

$$(29) \quad \psi[\theta^*(E), E] \geq \psi[\bar{\theta}(E), E]$$

for all sample points  $E$ . Let us consider the integral

$$(30) \quad I = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} \{W[\theta - \bar{\theta}(E)] - W[\theta - \theta^*(E)]\} p(E | \theta) d\theta dx_1 \dots dx_n.$$

Integrating (30) with respect to  $\theta$  we get

$$(31) \quad I = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} \{\psi[\bar{\theta}(E), E] - \psi[\theta^*(E) - E]\} dx_1 \dots dx_n.$$

Integrating (30) with respect to  $E$ , we get

$$(32) \quad I = \int_{-\infty}^{+\infty} [r(\theta) - r^*(\theta)] d\theta.$$

On account of (29) and (31) we have  $I \leq 0$ , hence

$$(33) \quad \int_{-\infty}^{+\infty} [r(\theta) - r^*(\theta)] d\theta \leq 0.$$

From (33) it follows that if  $r^*(\theta) \leq r(\theta)$  for every  $\theta$  then  $r^*(\theta) < r(\theta)$  can hold only for the points of a set of measure zero. In case  $r^*(\theta)$  is continuous, this means that  $r^*(\theta) \equiv r(\theta)$ . Hence if  $r^*(\theta)$  is continuous, then either  $r^*(\theta) \equiv r(\theta)$  or there exists at least one point  $\theta'$  such that  $r^*(\theta') > r(\theta')$ . The risk function  $r^*(\theta)$  is continuous if the estimate  $\theta^*(E)$  is uniformly continuous in the whole sample space. In fact, we have

$$r^*(\theta + t) = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} W[\theta + t - \theta^*(E)] p(x_1 - \theta - t, \dots, x_n - \theta - t) dx_1 \dots dx_n.$$

Making the transformation

$$y_i = x_i - t \quad (i = 1, \dots, n)$$

we get

$$r^*(\theta + t) = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} W[\theta + t - \theta^*(y_1 + t, \dots, y_n + t)] p(y_1 - \theta, \dots, y_n - \theta) dy_1 \dots dy_n.$$

Since  $W(u)$  and  $\theta^*(E)$  are uniformly continuous, from the latter equation we get easily

$$\lim_{t \rightarrow 0} r^*(\theta + t) = r^*(\theta)$$

that is to say  $r^*(\theta)$  is continuous. Considering only continuous estimates the admissibility of  $\bar{\theta}(E)$ , and therefore also Theorem 5, is proved. If  $\theta^*(E)$  is not

uniformly continuous we have only proved that if  $r^*(\theta) \leq r(\theta)$  for every  $\theta$ , then  $r^*(\theta) < r(\theta)$  can hold only in a set of measure zero. I should like to mention without proof that even if  $\theta^*(E)$  is not continuous,  $r^*(\theta) \leq r(\theta)$  implies  $r^*(\theta) \equiv r(\theta)$ .

An estimate  $\hat{\theta}(E)$  is called a maximum likelihood estimate if for any fixed  $E$   $p(E | \theta)$  becomes a maximum with respect to  $\theta$  for  $\theta = \hat{\theta}(E)$ .

**THEOREM 6.** *Consider the following conditions:*

VI. *There exists exactly one maximum likelihood estimate  $\hat{\theta}(E)$  with the following properties:*

a) *For any  $E$   $p(\hat{E} | \theta)$  is non-decreasing with increasing  $\theta$  for  $\theta < \hat{\theta}(E)$  and non-increasing with increasing  $\theta$  for  $\theta > \hat{\theta}(E)$ .*

b) *For any  $E$   $p(E | \theta)$  is a symmetric function of  $\theta$  about  $\hat{\theta}(E)$  that is to say, for for any real value  $\lambda$   $p[E | \hat{\theta}(E) - \lambda] = p[E | \hat{\theta}(E) + \lambda]$ .*

VII. *The value of the weight function depends only on the absolute value of the difference  $u = \theta - \bar{\theta}$  and  $\frac{dw(u)}{du}$  exists, is uniformly continuous and  $> 0$  for  $u > 0$ .*

*If the conditions I-V of Theorem 5 and the above condition VII are fulfilled, and if  $\hat{\theta}(E)$  is a maximum likelihood estimate satisfying VI, then  $\hat{\theta}(E)$  is a best estimate.*

Assume that the conditions I-V and VII are satisfied and that  $\hat{\theta}(E)$  is a maximum likelihood estimate satisfying VI. It is obvious that  $\hat{\theta}(E'') - \hat{\theta}(E') = \lambda$  for  $E' = (x_1, \dots, x_n)$  and  $E'' = (x_1 + \lambda, \dots, x_n + \lambda)$ . In order to prove Theorem 6, we have, according to Theorem 5, only to show that the integral in (28)

$$\psi(\bar{\theta}, E) = \int_{-\infty}^{+\infty} W(\theta - \bar{\theta})p(E | \theta) d\theta$$

becomes a minimum for  $\bar{\theta} = \hat{\theta}(E)$ . Denote  $\theta - \bar{\theta}$  by  $u$ . Since  $\frac{dW(u)}{du}$  is uniformly continuous, we have

$$\frac{\partial \psi(\bar{\theta}, E)}{\partial \bar{\theta}} = \int_{-\infty}^{+\infty} -\left[\frac{dW(u)}{du}\right] p(E | \theta) d\theta.$$

Since  $\frac{dW(u)}{du} = -\frac{dW(-u)}{du}$  we have

$$(34) \quad \frac{\partial \psi(\bar{\theta}, E)}{\partial \bar{\theta}} = \int_0^{\infty} \left[\frac{dW(u)}{du}\right] [p(E | \bar{\theta} - u) - p(E | \bar{\theta} + u)] du.$$

From condition VI it follows easily that for any fixed  $E$  and  $\bar{\theta}$  the function of  $u$  ( $0 \leq u \leq \infty$ )

$$p(E | \bar{\theta} - u) - p(E | \bar{\theta} + u)$$

does not change its sign and if  $\bar{\theta} \neq \hat{\theta}(E)$  there exists an interval  $J$  such that the above expression is unequal to zero for every point  $u$  of  $J$ . Hence on account of  $\frac{dW(u)}{du} > 0$  for  $u > 0$ , the integral in (34) vanishes only for  $\bar{\theta} = \hat{\theta}(E)$ . Since according to the condition V there exists a finite value  $\theta'$  such that  $\psi(\bar{\theta}, E)$  becomes a minimum for  $\bar{\theta} = \theta'$ ,  $\theta'$  must be equal to  $\hat{\theta}(E)$ . This proves Theorem 6.

The condition VI is seldom exactly fulfilled. But for large  $n$ , in the great majority of practical cases, VI will be fulfilled with good approximation and the best estimate approaches the maximum likelihood estimate with increasing  $n$ .

**10. Two examples.** As a first example we consider a normal distribution with the variance 1. The mean value  $\theta$  is unknown and we have to estimate it by means of a sample  $E = (x_1, \dots, x_n)$ . In this case

$$p(E|\theta) = \frac{1}{(2\pi)^{\frac{n}{2}}} e^{-\frac{1}{2}\sum(x_i-\theta)^2}.$$

It is obvious that for a very broad class of weight functions the conditions I-V of Theorem 5 are fulfilled. The maximum likelihood estimate  $\hat{\theta}(x_1, \dots, x_n) = \frac{x_1 + \dots + x_n}{n}$  satisfies the condition VI of Theorem 6. Hence if the weight function satisfies also the condition VII, then the best estimate of  $\theta$  is the maximum likelihood estimate  $\hat{\theta}(x_1, \dots, x_n) = \frac{x_1 + \dots + x_n}{n}$ .

Let us now consider a weight function defined as follows:

$$W(\theta, \bar{\theta}) = 2(\bar{\theta} - \theta) \quad \text{if } \bar{\theta} \geq \theta$$

and

$$W(\theta, \bar{\theta}) = \theta - \bar{\theta} \quad \text{if } \bar{\theta} < \theta.$$

Since for this weight function, the conditions I-V satisfied, according to Theorem 5 the best estimate of  $\theta$  is the value  $\bar{\theta}$  for which the integral

$$\int_{-\infty}^{+\infty} W(\theta, \bar{\theta}) e^{-\frac{1}{2}\sum(x_i-\theta)^2} d\theta = \int_{-\infty}^{\bar{\theta}} 2(\bar{\theta} - \theta) e^{-\frac{1}{2}\sum(x_i-\theta)^2} d\theta + \int_{\bar{\theta}}^{+\infty} (\theta - \bar{\theta}) e^{-\frac{1}{2}\sum(x_i-\theta)^2} d\theta$$

becomes a minimum. As an easy calculation shows, the estimate obtained in this way is not the arithmetic mean.

As a second example we consider the family of variates  $X(\theta)$  with the probability density  $f(x, \theta)$  defined as follows:

$$f(x, \theta) = 1 \quad \text{if } \theta - \frac{1}{2} \leq x \leq \theta + \frac{1}{2}$$

and

$$f(x, \theta) = 0 \quad \text{for all other values of } x.$$



If  $E = (x_1, \dots, x_n)$  denotes a sample point where  $x_1$  denotes the smallest and  $x_n$  denotes the greatest value in the sample, then

$$p(E | \theta) = \prod_{i=1}^n f(x_i, \theta) = 1 \quad \text{if } x_n - \frac{1}{2} \leq \theta \leq x_1 + \frac{1}{2}$$

and

$$p(E | \theta) = 0 \quad \text{for all other values of } \theta.$$

The classical method of maximum likelihood cannot be applied here, since  $p(E | \theta)$  is maximum for every value  $\theta$  for which  $x_n - \frac{1}{2} \leq \theta \leq x_1 + \frac{1}{2}$ . It is obvious that for a broad class of weight functions the conditions I-V are satisfied. The estimate  $\bar{\theta}(E) = \frac{x_1 + x_n}{2}$ , where  $x_1$  denotes the smallest and  $x_n$  the greatest value in the sample, satisfies the condition VI. Hence if the weight function satisfies also the condition VII, the best estimate of  $\theta$  is given by  $\bar{\theta}(E) = \frac{x_1 + x_n}{2}$ .

Let us now calculate the best estimate of  $\theta$  if the weight function is given as follows:

$$W(\theta, \bar{\theta}) = \theta - \bar{\theta} \quad \text{if } \bar{\theta} \leq \theta$$

and

$$W(\theta, \bar{\theta}) = 2(\bar{\theta} - \theta) \quad \text{if } \bar{\theta} > \theta.$$

In this case the conditions I-V are satisfied but not the condition VII. We have to calculate the integral  $\psi(\bar{\theta}, E)$  given in (28), which reduces in this case to

$$\begin{aligned} \psi(\bar{\theta}, E) &= \int_{x_n - \frac{1}{2}}^{x_1 + \frac{1}{2}} W(\theta, \bar{\theta}) d\theta = \int_{x_n - \frac{1}{2}}^{\bar{\theta}} 2(\bar{\theta} - \theta) d\theta + \int_{\bar{\theta}}^{x_1 + \frac{1}{2}} (\theta - \bar{\theta}) d\theta \\ &= 1.5\bar{\theta}^2 - [(x_1 + \frac{1}{2}) + 2(x_n - \frac{1}{2})]\bar{\theta} + \frac{1}{2}(x_1 + \frac{1}{2})^2 + (x_n - \frac{1}{2})^2. \end{aligned}$$

This expression becomes a minimum for

$$\bar{\theta} = \frac{x_1 + 2x_n - \frac{1}{2}}{3}.$$

Hence the best estimate of  $\theta$  is given by this expression.

**11. Miscellaneous remarks.** Assumptions 1-6 of paragraph 8 are sufficient but not necessary for the proof of the Theorems 1-3 (Theorems 4-6 have been deduced without Assumptions 1-6). They can be weakened in many respects. The assumption that the parameter space is bounded can be dropped if we impose certain conditions on the weight function  $W(\theta, \bar{\theta})$  and the probability density  $p(E | \theta)$ . It is certainly not necessary to assume that  $W(\theta, \bar{\theta})$  and  $p(E | \theta)$  are everywhere continuous. It is however doubtful whether Theorems 1-3 remain valid in the form in which they are stated, if we admit discon-

tinuities in a set of measure zero without imposing any other restrictions. Also Assumptions 4–6 can in all probability be essentially weakened.

I should like to mention that Assumption 4 is not as restrictive as it would appear. Let us make this clear in the case that the parameter space is a one-dimensional interval  $[a, b]$ . If we assume that  $W(\theta, \bar{\theta})$  is a polynomial of the second degree in  $\bar{\theta}$  and the coefficient of  $\bar{\theta}^2$  is positive for every  $\theta$ , and if  $p(E | \theta) > 0$  for every  $E$  and  $\theta$ , the Assumption 4 can easily be proved. In fact,

$$\psi(\bar{\theta}, E) = \int_a^b W(\theta, \bar{\theta}) p(E | \theta) df(\theta) = A(E) + B(E)\bar{\theta} + C(E)\bar{\theta}^2.$$

Since the coefficient of  $\bar{\theta}^2$  in  $W(\theta, \bar{\theta})$  is positive and since  $p(E | \theta) > 0$  for every  $E$  and  $\theta$ ,  $C(E) > 0$  for every  $E$  and for any arbitrary distribution  $f(\theta)$ . From this fact follows easily that for every  $E$  there exists a value  $\bar{\theta}(E)$  in the interval  $[a, b]$  such that

$$\psi[\bar{\theta}(E), E] < \psi(\bar{\theta}, E)$$

for every  $\bar{\theta}$  contained in  $[a, b]$  and unequal to  $\bar{\theta}(E)$ . Hence Assumption 4 is proved.

Let us consider a system  $S$  of subsets of the parameter space  $\Omega$  and the corresponding system  $H_s$  of hypotheses. The weight function  $W(\theta, \omega)$  is defined for all points  $\theta$  of  $\Omega$  and for all elements  $\omega$  of  $S$  and expresses the weight of the error committed by accepting  $H_\omega$  when  $\theta$  is true. If  $\theta$  is an element of  $\omega$  then  $W(\theta, \omega)$  is of course equal to zero. Let us assume that  $W(\theta, \omega)$  has the special form:  $W(\theta, \omega) = 1$  if  $\theta$  is not contained in  $\omega$ , and  $W(\theta, \omega) = 0$  if  $\theta$  is an element of  $\omega$ . It is obvious that in this case for any  $\theta$  the value of the risk function  $r(\theta)$  is equal to the probability of accepting a false hypothesis if  $\theta$  is the true parameter point. Because of this fact the theory developed here has close relation to the theory of confidence intervals. Let us first make this clear for the case when the parameter space is one dimensional, that is to say  $\theta$  is a real number.

In the theory of confidence intervals we estimate the unknown parameter  $\theta$  by an interval  $I(E)$  extending from  $\theta_1(E)$  to  $\theta_2(E)$  where  $\theta_1(E)$  and  $\theta_2(E)$  are certain functions of the sample point  $E$ . The interval  $I(E)$  is defined in such a way that the following probability statement holds: If we perform an experiment, the probability that we shall obtain a sample point  $E$  such that  $I(E)$  will cover the true parameter point  $\theta$ , is equal to a given constant  $\alpha$  (called confidence coefficient) and is independent of the value of  $\theta$ . Let us consider a certain example of such an inference with the confidence coefficient  $\alpha$  and denote by  $I(E)$  the interval corresponding to  $E$ . We define a system  $S$  of intervals as follows: An interval  $I$  is an element of  $S$  if and only if there exists a sample point  $E$  for which  $I(E) = I$ . Consider the corresponding system  $H_s$

of hypotheses and the weight function  $W(\theta, I)$  defined for all values  $\theta$  and all elements  $I$  of  $S$  as follows:

$$\begin{aligned} W(\theta, I) &= 0 & \text{if } \theta \text{ is a point of } I \\ W(\theta, I) &= 1 & \text{if } \theta \text{ is not contained in } I. \end{aligned}$$

Denote by  $M_s$  a best system of regions of acceptance relative to the weight function defined above. Denote by  $I'(E)$  the element of  $S$  which we accept according to  $M_s$  if  $E$  is the sample point. On account of the special form of the weight function, the risk is obviously equal to the probability of accepting a false interval. From the definition of the best system of regions it follows that for any  $\theta$  the probability that  $I'(E)$  will cover  $\theta$  is greater than or equal to  $\alpha$ . If the risk function is constant, that is to say, if the probability that  $I'(E)$  will cover the true parameter point  $\theta$  is independent of the value of  $\theta$ , then the intervals  $I'(E)$  are confidence intervals corresponding to a confidence coefficient  $\alpha' \geq \alpha$ .

Similar observations can be made if the parameter space is  $k$ -dimensional ( $k > 1$ ) that is to say,  $\theta$  is a system of  $k$  numbers  $\theta^{(1)}, \dots, \theta^{(k)}$ . An important case is that when we have to estimate only one of the components, say  $\theta^{(1)}$ , by an interval. As the investigations of W. Feller<sup>8</sup> have shown, confidence intervals in such cases do not exist always. That is to say, it is not always possible to determine  $I(E)$  such that the probability that  $I(E)$  will cover  $\theta^{(1)}$  is equal to a given constant  $\alpha$  independently of the values of  $\theta^{(1)}, \dots, \theta^{(k)}$ . It is of great interest to know under what conditions confidence intervals exist. I should like to mention that a further development of the theory given in paragraph 8 may contribute much to the solution of this problem. In order to make this clear, let us consider a system  $S_1$  of one-dimensional intervals. To each element  $I$  of  $S_1$  let there correspond the subset  $\omega$  of the  $k$ -dimensional parameter space  $\Omega$  consisting of all points  $\theta = (\theta^{(1)}, \dots, \theta^{(k)})$  for which  $\theta^{(1)}$  lies in  $I$ . Consider the system  $S$  of subsets  $\omega$  of  $\Omega$  corresponding to all elements of  $S_1$  and the system  $H_s$  of hypotheses corresponding to  $S$ . The weight function is to be chosen as follows:  $W(\theta, \omega) = 1$  if  $\theta$  is not an element of  $\omega$  and  $W(\theta, \omega) = 0$  if  $\theta$  is an element of  $\omega$ . Consider a best system  $M_s$  of regions of acceptance and the corresponding risk function  $r(\theta)$ . On account of the special definition of  $W(\theta, \omega)$ ,  $r(\theta)$  is equal to the probability of accepting a false hypothesis if  $\theta$  is the true parameter point. If the risk function  $r(\theta)$  is identically equal to a constant  $\alpha$ , we have confidence intervals corresponding to the confidence coefficient  $\alpha$ . In order to see under what conditions the risk function is constant, we have to consider an equivalent problem (see paragraph 7) where the system of hypotheses is the system of all simple hypotheses and the weight function  $W(\theta, \bar{\theta})$

<sup>8</sup> W. Feller, "Note on Regions Similar to the Sample Space," *Statistical Research Memoirs*, Vol. II, 1938.

is given according to formula (4). If  $W(\theta, \bar{\theta})$  satisfies Assumptions 1-6, the risk function of the best estimate is constant. As we have mentioned, Assumptions 1-6 can be weakened. In order to get valuable results concerning the problem of the existence of confidence intervals, we have to weaken especially Assumption 2. In fact  $W(\theta, \omega)$  takes only the values 1 and 0 and therefore  $W(\theta, \bar{\theta})$  cannot be continuous.

Finally I should like to mention that the most stringent test as defined by Robert W. B. Jackson<sup>9</sup> is contained as special case in our general definition of the best system of regions of acceptance. Jackson considers a discontinuous parameter space  $\Omega$ . Consider the problem of testing the hypothesis  $\theta = \theta_0$  where  $\theta_0$  denotes a point of  $\Omega$ . According to Jackson's definition we have the most stringent test if the critical region  $w_0$  satisfies the condition: the maximum of the numbers  $A$  and  $B$

$A = P(E \in w \mid \theta_0)$ ,  $B =$  least upper bound of  $P(E \in \bar{w} \mid \theta)$  formed for all  $\theta \neq \theta_0$ ,

becomes a minimum for  $w = w_0$ .  $\bar{w}$  denotes the region complementary to  $w$ . It is easy to see that Jackson's definition of the most stringent test coincides with our definition of the best system of regions of acceptance in the following special case:

- 1)  $\Omega$  is discontinuous
- 2)  $S$  consists only of two elements.
- 3) The weight function  $W(\theta, \omega)$  is equal to 1 if  $\theta$  is not contained in  $\omega$ .

COLUMBIA UNIVERSITY.

---

<sup>9</sup> Robert W. Jackson, "Testing Statistical Hypotheses," *Statistical Research Memoirs*, Vol. I, 1936.