

6 Mathematics of a Lady Tasting Tea By SIR RONALD A. FISHER

STATEMENT OF EXPERIMENT

A LADY declares that by tasting a cup of tea made with milk she can discriminate whether the milk or the tea infusion was first added to the cup: We will consider the problem of designing an experiment by means of which this assertion can be tested. For this purpose let us first lay down a simple form of experiment with a view to studying its limitations and its characteristics, both those which appear to be essential to the experimental method, when well developed, and those which are not essential but auxiliary .

Our experiment consists in mixing eight cups of tea, four in one way and four in the other, and presenting them to the subject for judgment in a random order. The subject has been told in advance of what the test will consist, namely that she will be asked to taste eight cups, that these shall be four of each kind, and that they shall be presented to her in a random order, that is in an order not determined arbitrarily by human choice, but by the actual manipulation of the physical apparatus used in games of chance, cards, dice, roulettes, etc., or, more expeditiously, from a published collection of random sampling numbers purporting to give the actual results of such manipulation. Her task is to divide the 8 cups into two sets of 4, agreeing, if possible, with the treatments received.

INTERPRETATION AND ITS REASONED BASIS

In considering the appropriateness of any proposed experimental design, it is always needful to forecast all possible results of the experiment, and to have decided without ambiguity what interpretation shall be placed upon each one of them. Further, we must know by what argument this interpretation is to be sustained . In the present instance we may argue as follows. There are 70 ways of choosing a group of 4 objects out of 8. This may be demonstrated by an argument familiar to students of "permutations and combinations," namely, that if we were to choose the 4 objects in succession we should have successively 8, 7, 6, 5 objects to choose from, and could make our succession of choices in $8 \times 7 \times 6 \times 5$, or 1680 ways. But in doing this we have not only chosen every possible set of 4, but every possible set in every possible order; and since 4 objects can be arranged in order in $4 \times 3 \times 2 \times 1$, or 24 ways, we may find the number of possible choices by dividing 1680 by 24. The result, 70, is essential to our interpretation of the experiment. At best the subject can judge rightly with every cup and, knowing that 4 are of each kind, this amounts to claim, out of the 70 sets of 4 which might be chosen, that particular one which is correct. A subject without any faculty of discrimination would in fact divide the 8 cups correctly into two sets of 4 in one trial out of 70, or, more properly, with a frequency which would approach 1 in 70 more and more nearly the more often the test were repeated. Evidently this frequency, with which unfailling success would be achieved by a person lacking altogether the faculty under test, is calculable from the number of cups used. The odds could be made much higher by enlarging the experiment, while, if the experiment were much smaller even the greatest possible success would give odds so low that the result might, with considerable probability, be ascribed to chance.

THE TEST OF SIGNIFICANCE

It is open to the experimenter to be more or less exacting in respect of the smallness of the probability he would require before he would be willing to admit that his observations have demonstrated a positive result. It is obvious that an experiment would be useless of which no possible result would satisfy him. Thus, if he wishes to ignore results having probabilities as high as 1 in 20 — the probabilities being of course reckoned from the hypothesis that the phenomenon to be demonstrated is in fact absent — then it would be useless for him to experiment with only 3 cups of tea of each kind. For 3 objects can be chosen out of 6 in only 20 ways and therefore complete success in the test would be achieved without sensory discrimination, i.e., by "pure chance," in an average of 5 trials out of 100. It is usual and convenient for experimenters to take 5 percent as a standard level of significance, in the sense that they are prepared to ignore all results which fail to reach this standard, and, by this means, to eliminate from further discussion the greater part of the fluctuations which chance causes have introduced into their experimental results. No such selection can eliminate the whole of the possible effects of chance coincidence, and if we accept this convenient convention, and agree that an event which would occur by chance only once in 70 trials is decidedly significant," in the statistical sense, we thereby admit that no isolated experiment, however significant in itself, can suffice for the experimental demonstration of any natural phenomenon; for the "one chance in a million" will undoubtedly occur, with no less and no more than its appropriate frequency, however surprised we may be that it should occur to us. In order to assert that a natural phenomenon is experimentally demonstrable we need, not an isolated record, but a reliable method of procedure. In relation to the test of significance, we may say that a phenomenon is experimentally demonstrable when we know how to conduct an experiment which will rarely fail to give us a statistically significant result.

Returning to the possible results of the psycho-physical experiment, having decided that if every cup were rightly classified a significant positive result would be recorded, or, in other words, that we should admit that the lady had made good her claim, what should be our conclusion if, for each kind of cup, her judgments are 3 right and 1 wrong? We may take it, in the present discussion, that any error in one set of judgments will be compensated by an error in the other, since it is known to the subject that there are 4 cups of each kind. In enumerating the number of ways of choosing 4 things out of 8, such that 3 are right and 1 wrong, we may note that the 3 right may be chosen, out of the 4 available, in 4 ways and, independently of this choice, that the 1 wrong may be chosen, out of the 4 available, also in 4 ways. So that in all we could make a selection of the kind supposed in 16 different ways. A similar argument shows that, in each kind of judgment, 2 may be right and 2 wrong in 36 ways, 1 right and 3 wrong in 16 ways and none right and 4 wrong in 1 way only. It should be noted that the frequencies of these five possible results of the experiment make up together, as it is obvious they should, the 70 cases out of 70.

It is obvious too, that 3 successes to 1 failure, although showing a bias, or deviation in the right direction, could not be judged as statistically significant evidence of a real sensory discrimination. For its frequency of chance occurrence is 16 in 70, or more than 20 per cent. Moreover, it is not the best possible result, and in judging of its significance we must take account not only of its own frequency, but also of the frequency for any better result. In the present instance "3 right and 1 wrong" occurs 16 times, and "4 right" occurs once in 70 trials, making 17 cases out of 70 as good as or better than that observed. The reason for including cases better than that observed becomes obvious on considering what our conclusions would have been had the case of 3 right and 1 wrong only 1 chance, and the case of 4 right 16 chances of occurrence out of 70. The rare case of 3 right and 1 wrong could not be judged significant merely because it was rare, seeing that a higher degree of success would frequently have been scored by mere chance.

THE NULL HYPOTHESIS

Our examination of the possible results of the experiment has therefore led us to a statistical test of significance, by which these results are divided into two classes with opposed interpretations. Tests of significance are of many different kinds, which need not be considered here. Here we are only concerned with the fact that the easy calculation in permutations which we encountered, and which gave us our test of significance, stands for something present in every possible experimental arrangement; or, at least, for something required in its interpretation. The two classes of results which are distinguished by our test of significance are, on the one

hand, those which show a significant discrepancy from a certain hypothesis; namely, in this case, the hypothesis that the judgments given are in no way influenced by the order in which the ingredients have been added; and on the other hand, results which show no significant discrepancy from this hypothesis. This hypothesis, which may or may not be impugned by the result of an experiment, is again characteristic of all experimentation. Much confusion would often be avoided if it were explicitly formulated when the experiment is designed. In relation to any experiment we may speak of this hypothesis as the "null hypothesis," and it should be noted that the null hypothesis is never proved or established, but is possibly disproved, in the course of experimentation. Every experiment may be said to exist only in order to give the facts a chance of disproving the null hypothesis.

It might be argued that if an experiment can disprove the hypothesis that the subject possesses no sensory discrimination between two different sorts of objects, it must therefore be able to prove the opposite hypothesis, that she can make some such discrimination. But this last hypothesis, however reasonable or true it may be, is ineligible, as a null hypothesis to be tested by experiment, because it is inexact. If it were asserted that the subject would never be wrong in her judgments we should again have an exact hypothesis, and it is easy to see that this hypothesis could be disproved by a single failure, but could never be proved by any finite amount of experimentation. It is evident that the null hypothesis must be exact, that is free from vagueness and ambiguity, because it must supply the basis of the "problem of distribution," of which the test of significance is the solution. A null hypothesis may, indeed, contain arbitrary elements, and in more complicated cases often does so: as, for example, if it should assert that the death-rates of two groups of animals are equal without specifying what these death-rates usually are. In such cases it is evidently the equality rather than any particular values of the death-rates that the experiment is designed to test, and possibly to disprove.

In cases involving statistical "estimation" these ideas may be extended to the simultaneous considerations of a series of hypothetical possibilities. The notion of an error of the so-called "second kind," due to accepting the null hypothesis "when it is false" may then be given a meaning in reference to the quantity to be estimated. It has no meaning with respect to simple tests of significance, in which the only available expectations are those which flow from the null hypothesis being true.

RANDOMIZATION; THE PHYSICAL BASIS OF THE VALIDITY OF THE TEST

We have spoken of the experiment as testing a certain null hypothesis, namely, in this case, that the subject possesses no sensory discrimination whatever of the kind claimed; we have, too, assigned as appropriate to this hypothesis a certain frequency distribution of occurrences, based on the equal frequency of the 70 possible ways of assigning 8 objects to two classes of 4 each; in other words, the frequency distribution appropriate to a classification by pure chance. We have now to examine the physical conditions of the experimental technique needed to justify the assumption that, if discrimination of the kind under test is absent, the result of the experiment will be wholly governed by the laws of chance. It is easy to see that it might well be otherwise. If all those cups made with the milk first had sugar added, while those made with the tea first had none, a very obvious difference in flavour would have been introduced which might well ensure that all those made with sugar should be classed alike. These groups might either be classified all right or all wrong, but in such a case the frequency of the critical event in which all cups are classified correctly would not be 1 in 70, but 35 in 70 trials, and the test of significance would be wholly vitiated. Errors equivalent in principle to this are very frequently incorporated in otherwise well-designed experiments.

It is no sufficient remedy to insist that "all the cups must be exactly alike" in every respect except that to be tested. For this is a totally impossible requirement in our example, and equally in all other forms of experimentation. In practice it is probable that the cups will differ perceptibly in the thickness or smoothness of their material, that the quantities of milk added to the different cups will not be exactly equal, that the strength of the infusion of tea may change between pouring the first and the last cup, and that the temperature also at which the tea is tasted will change during the course of the experiment. These are only examples of the differences probably present; it would be impossible to present an exhaustive list of such possible differences appropriate to any one kind of experiment, because the uncontrolled causes which may influence the result are always strictly innumerable. When any such cause is named, it is usually perceived that, by increased labour and expense, it could be largely eliminated. Too frequently it is assumed that such refinements constitute improvements to the experiment. Our view, which will be much more fully exemplified in later sections, is that it is an essential

characteristic of experimentation that it is carried out with limited resources, and an essential part of the subject of experimental design to ascertain how these should be best applied: or, in particular, to which causes of disturbance care should be given, and which *ought* to be deliberately ignored. To ascertain, too, for those which are not to be ignored, to what *extent* it is worth while to take the trouble to diminish their magnitude. For our present purpose, however, it is only necessary to recognize that, whatever degree of care and experimental skill is expended in equalizing the conditions, other than the one under test, which are liable to affect the result, this equalization must always be to a greater or less extent incomplete, and in many important practical cases will certainly be grossly defective. We are concerned, therefore, that this inequality, whether it be great or small, shall not impugn the exactitude of the frequency distribution, on the basis of which the result of the experiment is to be appraised .

THE EFFECTIVENESS OF RANDOMIZATION

The element in the experimental procedure which contains the essential safeguard is that the two modifications of the test beverage are to be prepared "in random order." This, in fact, is the only point in the experimental procedure in which the laws of chance, which are to be in exclusive control of our frequency distribution, have been explicitly introduced. The phrase "random order" itself, however, must be regarded as an incomplete instruction, standing as a kind of shorthand symbol for the full procedure of randomization, by which the validity of the test of significance may be guaranteed against corruption by the causes of disturbance which have not been eliminated. To demonstrate that, with satisfactory randomization, its validity is, indeed, wholly unimpaired, let us imagine all causes of disturbances--the strength of the infusion, the quantity of milk, the temperature at which it is tasted, etc .--to be predetermined for each cup; then since these, on the null hypothesis, are the only causes influencing classification, we may say that the probabilities of each of the 70 possible choices or classifications which the subject can make are also predetermined. If, now, after the disturbing causes are fixed, we assign, strictly at random 4 out of the 8 cups to each of our experimental treatments, then every set of 4, whatever its probability of being so classified, will certainly have a probability of exactly 1 in 70 of being the 4, for example, to which the milk is added first. However important the causes of disturbance may be, even if they were to make it certain that one particular set of 4 should receive this classification, the probability that the 4 so classified and the 4 which ought to have been so classified should be the same, must be rigorously in accordance with our test of significance.

It is apparent, therefore, that the random choice of the objects to be treated in different ways would be a complete guarantee of the validity of the test of significance, if these treatments were the last in time of the stages in the physical history of the objects which might affect their experimental reaction. The circumstance that the experimental treatments cannot always be applied last, and may come relatively early in their history, causes no practical inconvenience; for subsequent causes of differentiation, if under the experimenter's control, as, for example, the choice of different pipettes to be used with different flasks, can either be predetermined before the treatments have been randomized, or, if this has not been done, can be randomized on their own account ; and other causes of differentiation will be either (a) consequences of differences already randomized, or (b) natural consequences of the differences in treatment to be treated, of which on the null hypothesis there will be none, by definition, or (c) effects supervening by chance independently from the treatments applied. Apart, therefore, from the avoidable error of the experimenter himself introducing with his test treatments, or subsequently, other differences in treatment, the effects of which the experiment is not intended to study, it may be said that the simple precaution of randomization will suffice to guarantee the validity of the test of significance, by which the result of the experiment is to be judged.

THE SENSITIVENESS OF AN EXPERIMENT. EFFECTS OF ENLARGEMENT AND REPETITION

A probable objection, which the subject might well make to the experiment so far described, is that only if every cup is classified correctly will she be judged successful. A single mistake will reduce her performance below the level of significance. Her claim, however, might be, not that she could draw the distinction with invariable certainty, but that, though sometimes mistaken, she would be right more often than not; and that the experiment should be enlarged sufficient-

ly, or repeated sufficiently often, for her to be able to demonstrate the predominance of correct classifications in spite of occasional errors.

An extension of the calculation upon which the test of significance was based shows that an experiment with 12 cups, six of each kind, gives, on the null hypothesis, 1 chance in 924 for complete success, and 36 chances for 5 of each kind classified right and 1 wrong. As 37 is less than a twentieth of 924, such a test could be counted as significant, although a pair of cups have been wrongly classified; and it is easy to verify that, using larger numbers still, a significant result could be obtained with a still higher proportion of errors. By increasing the size of the experiment, we can render it more sensitive, meaning by this that it will allow of the detection of a lower degree of sensory discrimination, or, in other words, of a quantitatively smaller departure from the null hypothesis. Since in every case the experiment is capable of disproving, but never of proving this hypothesis, we may say that the value of the experiment is increased whenever it permits the null hypothesis to be more readily disproved.

The same result could be achieved by repeating the experiment, as originally designed, upon a number of different occasions, counting as a success all those occasions on which 8 cups are correctly classified. The chance of success on each occasion being 1 in 70, a simple application of the theory of probability shows that 2 or more successes in 10 trials would occur, by chance, with a frequency below the standard chosen for testing significance; so that the sensory discrimination would be demonstrated, although, in 8 attempts out of 10, the subject made one or more mistakes. This procedure may be regarded as merely a second way of enlarging the experiment and, thereby, increasing its sensitiveness, since in our final calculation we take account of the aggregate of the entire series of results, whether successful or unsuccessful. It would clearly be illegitimate, and would rob our calculation of its basis, if the unsuccessful results were not all brought into the account.

QUALITATIVE METHODS OF INCREASING SENSITIVENESS

Instead of enlarging the experiment we may attempt to increase its sensitiveness by qualitative improvements; and these are, generally speaking of two kinds: (a) the reorganization of its structure, and (b) refinements of technique. To illustrate a change of structure we might consider that, instead of fixing in advance that 4 cups should be of each kind, determining by a random process how the subdivision should be effected, we might have allowed the treatment of each cup to be determined independently by chance, as by the toss of a coin, so that each treatment has an equal chance of being chosen. The chance of classifying correctly 8 cups randomized in this way, without the aid of sensory discrimination, is 1 in 2^8 , or 1 in 256 chances, and there are only 8 chances of classifying 7 right and 1 wrong; consequently the sensitiveness of the experiment has been increased, while still using only 8 cups, and it is possible to score significant success, even if one is classified wrongly. In many types of experiment, therefore, the suggested change in structure would be evidently advantageous. For the special requirements of a psycho-physical experiment, however, we should probably prefer to forego this advantage, since it would occasionally occur that all the cups would be treated alike, and this, besides bewildering the subject by an unexpected occurrence, would deny her the real advantage of judging by comparison.

Another possible alteration to the structure of the experiment, which would, however, decrease its sensitiveness, would be to present determined, but unequal, numbers of the two treatments. Thus we might arrange that 5 cups should be of the one kind and 3 of the other, choosing them properly by chance, and informing the subject how many of each to expect. But since the number of ways of choosing 3 things out of 8 is only 56, there is now, on the null hypothesis, a probability of a completely correct classification of 1 in 56. It appears in fact that we cannot by these means do better than by presenting the two treatments in equal numbers, and the choice of this equality is now taken to be justified by its giving to the experiment its maximal sensitiveness.

With respect to the refinements of technique, we have seen above that these contribute nothing to the validity of the experiment, and of the test of significance by which we determine its result. They may, however, be important, and even essential, in permitting the phenomenon under test to manifest itself. Though the test of significance remains valid, it may be that without special precautions even a definite sensory discrimination would have little chance of scoring a significant success. If some cups were made with India and some with China tea, even though the treatments were properly randomized, the subject might not be able to discriminate

the relatively small difference in flavour under investigation, when it was confused with the greater differences between leaves of different origin. Obviously, a similar difficulty could be introduced by using in some cups raw milk and in others boiled, or even condensed milk, or by adding sugar in unequal quantities. The subject has a right to claim, and it is in the interests of the sensitiveness of the experiment, that gross differences of these kinds should be excluded, and that the cups should not as far as *possible*, but as far as is practically convenient, be made alike in all respects except that under test.

How far such experimental refinements should be carried is entirely a matter of judgment, based on experience. The validity of the experiment is not affected by them. Their sole purpose is to increase its sensitiveness, and this object can usually be achieved in many other ways, and particularly by increasing the size of the experiment. If, therefore, it is decided that the sensitiveness of the experiment should be increased, the experimenter has the choice between different methods of obtaining equivalent results; and will be wise to choose whichever method is easiest to him, irrespective of the fact that previous experimenters may have tried, and recommended as very important, or even essential, various ingenious and troublesome precautions.