

## STATISTICAL SIGNIFICANCE IN PSYCHOLOGICAL RESEARCH

DAVID T. LYKKEN

*University of Minnesota*

Most theories in the areas of personality, clinical, and social psychology predict no more than the direction of a correlation, group difference, or treatment effect. Since the null hypothesis is never strictly true, such predictions have about a 50-50 chance of being confirmed by experiment when the theory in question is false, since the statistical significance of the result is a function of the sample size. Confirmation of a single directional prediction should usually add little to one's confidence in the theory being tested. Most theories should be tested by multiple corroboration and most empirical generalizations by constructive replication. Statistical significance is perhaps the least important attribute of a good experiment; it is never a sufficient condition for claiming that a theory has been usefully corroborated, that a meaningful empirical fact has been established, or that an experimental report ought to be published.

In a recent journal article Sapolsky (1964) developed the following substantive theory: Some psychiatric patients entertain an unconscious belief in the "cloacal theory of birth" which involves the notions of oral impregnation and anal parturition. Such patients should be inclined to manifest eating disorders: compulsive eating in the case of those who wish to get pregnant and anorexia in those who do not. Such patients should also be inclined to see cloacal animals, such as frogs, on the Rorschach. This reasoning led Sapolsky to predict that Rorschach frog responders show a higher incidence of eating disorders than patients not giving frog responses. A test of this hypothesis in a psychiatric hospital showed that 19 of 31 frog responders had eating disorders indicated in their charts, compared to only 5 of the 31 control patients. A highly significant chi-square was obtained.

It will be an expository convenience to analyze Sapolsky's article in considerable detail for purposes of illustrating the methodological issues which are the real subject of this paper. My intent is not to criticize a particular author but rather to examine a kind of epistemic confusion which seems to be endemic in psychology, especially, but by no means exclusively, in its "softer" precincts. One would like to demonstrate this generality with multiple examples. Having just combed the latest issues of four well-known journals in the clinical and personality areas, I could undertake to identify several papers in each issue wherein, because they were able to

ject a directional null hypothesis at some high level of significance, the authors claimed to have usefully corroborated some rather general theory or to have demonstrated some important empirical relationship. To substantiate that these claims are overstated and that much of this research has not yet earned the right to the reader's overburdened attentions would require a lengthy analysis of each paper. Such profligacy of space would ill become an essay one aim of which is to restrain the swelling volume of the psychological literature. Therefore, with apologies to Sapolsky for subjecting this one paper to such heavy handed scrutiny, let us proceed with the analysis.

Since I regarded the prior probability of Sapolsky's theory (that frog responders unconsciously believe in impregnation per os) to be nugatory and its likelihood unenhanced by the experimental findings, I undertook to check my own reaction against that of 20 colleagues, most of them clinicians, by means of a formal questionnaire. The 20 estimates of the prior probability of Sapolsky's theory, which these psychologists made before being informed of his experimental results, ranged from  $10^{-6}$  to 0.13 with a median value of 0.01, which can be interpreted to mean, roughly, "I don't believe it." Since the prior probability of many important scientific theories is considered to be vanishingly small when they are first propounded, this result provides no basis for alarm. However, after being given a fair summary of Sapolsky's experimental findings, which "corroborate" the

theory by confirming the operational hypothesis derived from it with high statistical significance, these same psychologists attached posterior probabilities to the theory which ranged from  $10^{-5}$  to 0.14, with the median unchanged at 0.01. I interpret this consensus to mean, roughly, "I still don't believe it." This finding, I submit, *is* alarming because it signifies a sharp difference of opinion between, for example, the consulting editors of the journal and a substantial segment of its readership, a difference on the very fundamental question of what constitutes good (i.e., publishable) clinical research.

The thesis of the present paper is that Sapolsky and the editors were in fact following, with reasonable consistency, our traditional rules for evaluating psychological research, but that, as the Sapolsky paper exemplifies, at least two of these rules should be reconsidered. One of the rules examined here asserts roughly the following: "When a prediction or hypothesis derived from a theory is confirmed by experiment, a non-trivial increment in one's confidence in that theory should result, especially when one's prior confidence is low." Clearly, my 20 colleagues were violating this rule here since their confidence in the frog responder-cloacal birth theory was not, on the average, increased by the contemplation of Sapolsky's highly significant chi-square. From their comments it seems that they found it too hard to accept that a belief in oral impregnation could lead to frog responding merely because the frog has a cloacum. (One must, after all, admit that few patients know what a cloacum is or that a frog has one and that those few who do know probably will also know that the frog's eggs are both fertilized and hatched externally so neither oral impregnation nor anal birth are in any way involved. Hence, *neither* the average patient *nor* the biologically sophisticated patient should logically be expected to employ the frog as a symbol for an unconscious belief in oral conception.) My colleagues, on the contrary, found it relatively easy to believe that the observed association between frog responding and eating problems might be due to some other cause entirely (e.g., both symptoms are immature or regressive in character; the frog, with its

disproportionately large mouth and voice may well constitute a common orality totem and hence be associated with problems in the oral sphere; "squeamish" people might tend both to see frogs and to have eating problems; and so on.)

Assuming that this first rule *is* wrong in this instance, perhaps it could be amended to allow one to make exceptions in cases resembling this illustration. For example, one could add the codicil: "This rule may be ignored whenever one considers the theory in question to be overly improbable or whenever one can think of alternative explanations for the experimental results." But surely such an amendment would not do. ESP, for example, could never become scientifically respectable if the first exception were allowed, and one consequence of the second would be that the importance attached to one's findings would always be inversely related to the ingenuity of one's readers. The burden of the present argument is that this rule is wrong not only in a few exceptional instances *but as it is routinely applied to the majority of experimental reports in the psychological literature.*

#### CORROBORATING THEORIES BY EXPERIMENTAL CONFIRMATION OF THEORETICAL PREDICTIONS<sup>1</sup>

Most psychological experiments are of three kinds: (a) studies of the effect of some treatment on some output variables, which can be regarded as a special case of (b) studies of the difference between two or more groups of individuals with respect to some variable, which in turn are a special case of (c) the study of the relationship or correlation between two or more variables within some specified population. Using the bivariate correlation design as paradigmatic, then, one notes first that the strict null hypothesis must always be assumed to be false (this idea is not new and has recently been illuminated by Baken, 1966). Unless one of the variables is wholly unreliable so that the values obtained are strictly random, it would be foolish to suppose that the correlation between any two

<sup>1</sup> Much of the argument in this section is based upon ideas developed in certain unpublished memoranda by P. E. Meehl (personal communication, 1963) and in a recent article (Meehl, 1967).

variables is identically equal to 0.0000 . . . (or that the effect of some treatment or the difference between two groups is exactly *zero*). The molar dependent variables employed in psychological research are extremely complicated in the sense that the measured value of such a variable tends to be affected by the interaction of a vast number of factors, both in the present situation and in the history of the subject organism. It is exceedingly unlikely that any two such variables will not share at least some of these factors and equally unlikely that their effects will exactly cancel one another out.

It might be argued that the more complex the variables the smaller their average correlation ought to be since a larger pool of common factors allows more chance for mutual cancellation of effects in obedience to the Law of Large Numbers. However, one knows of a number of unusually potent and pervasive factors which operate to unbalance such convenient symmetries and to produce correlations large enough to rival the effects of whatever causal factors the experimenter may have had in mind. Thus, we know that (a) "good" psychological and physical variables tend to be positively correlated; (b) experimenters, without deliberate intention, can somehow subtly bias their findings in the expected direction (Rosenthal, 1963); (c) the effects of common method are often as strong as or stronger than those produced by the actual variables of interest (e.g., in a large and careful study of the factorial structure of adjustment to stress among officer candidates, Holtzman & Bitterman, 1956, found that their 101 original variables contained five main common factors representing, respectively, their rating scales, their perceptual-motor tests, the McKinney Reporting Test, their GSR variables, and the MMPI); (d) transitory state variables such as the subject's anxiety level, fatigue, or his desire to please, may broadly affect all measures obtained in a single experimental session.

This average shared variance of "unrelated" variables can be thought of as a kind of ambient noise level characteristic of the domain. It would be interesting to obtain empirical estimates of this quantity in our field to serve as a kind of Plimsoll mark

against which to compare obtained relationships predicted by some theory under test. If, as I think, it is not unreasonable to suppose that "unrelated" molar psychological variables share on the average about 4% to 5% of common variance, then the expected correlation between any such variables would be about .20 in absolute value and the expected difference between any two groups on some such variable would be nearly 0.5 standard deviation units. (Note that these estimates assume zero measurement error. One can better explain the near-zero correlations often observed in psychological research in terms of unreliability of measures than in terms of the assumption that the true scores are in fact unrelated.)

Suppose now that an investigator predicts that two variables are positively correlated. Since we expect the null hypothesis to be false, we expect his prediction to be confirmed by experiment with a probability of very nearly 0.5; by using a large enough sample, moreover, he can achieve any desired level of statistical significance for this result. If the ambient noise level for his domain is represented by correlations averaging, say, .20 in absolute value, then his chances of finding a statistically significant confirmation of his prediction with a reasonable sample size will be quite high (e.g., about 1 in 4 for  $N = 100$ ) even if there is no truth whatever to the theory on which the prediction was based. Since most theoretical predictions in psychology, especially in the areas of clinical and personality research, specify no more than the direction of a correlation, difference or treatment effect, we must accept the harsh conclusion that a single experimental finding of this usual kind (confirming a directional prediction), no matter how great its statistical significance, will seldom represent a large enough increment of corroboration for the theory from which it was derived to merit very serious scientific attention. (In the natural sciences, this problem is far less severe for two reasons: (a) theories are powerful enough to generate point predictions or at least predictions of some narrow range within which the dependent variable is expected to lie; and (b) in these sciences, the degree of experimental control and the relative sim-

plicity of the variables studied are such that the ambient noise level represented by unexplained and unexpected correlations, differences, and treatment effects is often vanishingly small.)

#### THE SIGNIFICANCE OF LARGE CORRELATIONS

It might be argued that, even where only a weak directional prediction is made, the obtaining of a result which is not only statistically significant but large in absolute value should constitute a stronger corroboration of the theory. For example, although Sapolsky predicted only that frog responding and eating disorders would be positively related, the fourfold point correlation (phi coefficient) between these variables in his sample was about .46, surely much larger than the average relationship expected between random pairs of molar variables on the premise that "everything is related to everything else." Does not such a large effect therefore provide stronger corroboration for the theory in question?

One difficulty with this reasonable sounding doctrine is that, in the complex sort of research considered here, *really large* effects, differences, or relationships are not usually to be expected and, when found, may even argue *against* the theory being tested. To illustrate this, let us take Sapolsky's theory seriously and, by making reasonable guesses concerning the unknown base rates involved, attempt to estimate the actual size of the relationship between frog responding and eating disorders which the theory should lead us to expect. Sapolsky found that 16% of his control sample showed eating disorders; let us take this value as the base rate for this symptom among patients who do not hold the cloacal theory of birth. Perhaps we can assume that all patients who do hold this theory will give frog responses but surely not all of these will show eating disorders (any more than will all patients who believe in vaginal conception be inclined to show coital or urinary disturbances); it seems a reasonable assumption that no more than 50% of the believers in oral conception will therefore manifest eating problems. Similarly, we can hardly suppose that the frog response *always* implies an unconscious belief in the cloacal

theory; surely this response can come to be emitted now and then for other reasons. Even with the greatest sympathy for Sapolsky's point of view, we could hardly expect more than, say, 50% of frog responders to believe in oral impregnation. Therefore, we might reasonably predict that 16 of 100 nonresponders would show eating disorders in a test of this theory, 50 of 100 frog responders would hold the cloacal theory and half of these show eating disorders, while 16% or 8 of the remaining 50 frog responders will show eating problems too, giving a total of 33 eating disorders among the 100 frog responders. Such a finding would produce a significant chi-square but the actual degree of relationship as indexed by the phi coefficient would be only about .20. In other words, if one considers the supplementary assumptions which would be required to make a theory compatible with the actual results obtained, it becomes apparent that the finding of a really strong association may actually embarrass the theory rather than support it (e.g., Sapolsky's finding of 61% eating disorders among his frog responders is *significantly larger* ( $p < .01$ ) than the 33% generously estimated by the reasoning above).

#### MULTIPLE CORROBORATION

In the social, clinical, and personality areas especially, we must expect that the size of the correlations, differences, or effects which might reasonably be predicted from our theories will typically not be very large relative to the ambient noise level of correlations and effects due solely to the "all-of-a-pieceness of things." The conclusion seems inescapable that the only really satisfactory solution to the problem of corroborating such theories is that of *multiple corroboration*, the derivation and testing of a number of separate, quasi-independent predictions. Since the prior probability of such a multiple corroboration may be on the order of  $(0.5)^n$ , where  $n$  is the number of independent<sup>2</sup> predictions experimentally confirmed, a theory of any useful degree of predictive richness should in principle allow

<sup>2</sup> Tests of predictions from the same theory are seldom strictly independent since they often share some of the same supplementary assumptions, are made at the same time on the same sample, and so on.

for sufficient empirical confirmation through multiple corroboration to compel the respect of the most critical reader or editor.

#### THE RELATION OF EXPERIMENTAL FINDINGS TO EMPIRICAL FACTS

We turn now to the examination of a second popular rule for the evaluation of psychological research, which states roughly that "When no obvious errors of sampling or experimental method are apparent, one's confidence in the general proposition being tested (e.g., Variables A and B are positively correlated in Population C) should be proportional to the degree of statistical significance obtained." We are following this rule when we say, "Theory aside, Sapolsky has at least demonstrated an empirical fact, namely, that frog responders have more eating disturbances than patients in general." This conclusion means, of course, that in the light of Sapolsky's highly significant findings we should be willing to give very generous odds that any other competent investigator (at another hospital, administering the Rorschach in his own way, and determining the presence of eating problems in whatever manner seems reasonable and convenient for him) will also find a substantial positive relationship between these two variables.

Let us be more specific. Given Sapolsky's fourfold table showing 19 of 31 frog responders to have eating disorders (61%), it can be shown by chi-square that we should have 99% confidence that the true population value lies between 13/31 and 25/31 (between 42% and 81%). With 99% confidence that the population value is at least 13 in 31, we should have  $.99(99) = 98\%$  confidence that a new sample from that population should produce at least 6 eating disorders among each 31 frog responders, assuming that 5 of each 31 nonresponders show eating problems also as Sapolsky reported. That is, we should be willing to bet \$98 against only \$2 that a replication of this experiment will show *at least as many* eating disorders among frog responders as among nonresponders. The reader may decide for himself whether his faith in the "empirical fact" demonstrated by this experiment can meet the test of this gambler's challenge.

#### THREE KINDS OF REPLICATION

If, as suggested above, "demonstrating an empirical fact" must involve a claim of confidence in the replicability of one's findings, then to clearly understand the relation of statistical significance to the probability of a "successful" replication it will be helpful to distinguish between three rather different methods of replicating or cross-validating an experiment. *Literal replication*, of course, would involve exact duplication of the first investigator's sampling procedure, experimental conditions, measuring techniques, and methods of analysis; asking the original investigator to simply run more subjects would perhaps be about as close as we could come to attaining literal replication and even this, in psychological research, might often not be close enough. In the case of *operational replication*, on the other hand, one strives to duplicate exactly just the sampling and experimental procedures given in the first author's report of his research. The purpose of operational replication is to test whether the investigator's "experimental recipe"—the conditions and procedures he considered salient enough to be listed in the "Methods" section of his report—will in other hands produce the results that he obtained. For example, replication of the "Clever Hans" experiment revealed that the apparent ability of that remarkable horse to add numbers had been due to an uncontrolled and unsuspected factor (the presence of the horse's trainer within his field of view). This factor, not being specified in the "methods recipe" for the result, was omitted in the replication which for that reason failed. Operational replication would be facilitated if investigators would accept more responsibility for specifying what they believe to be the minimum essential conditions and controls for producing their results. Psychologists tend to be inconsistently prolix in describing their experimental methods; thus, Sapolsky tabulates the age, sex, and diagnosis for each of his 62 subjects. Does he mean to imply that the experiment will not work if these details are changed?—surely not, but then why describe them?

In the quite different process of *constructive replication*, one deliberately avoids imitation

of the first author's methods. To obtain an ideal constructive replication, one would provide a competent investigator with *nothing more than* a clear statement of the empirical "fact" which the first author would claim to have established—for example, "psychiatric patients who give frog responses on the Rorschach have a greater tendency toward eating disorders than do patients in general"—and then let the replicator formulate his own methods of sampling, measurement, and data analysis. One must keep in mind that the data, the specific results of a particular experiment, are only seldom of any real interest in themselves. The "empirical facts" which we value so highly consist usually of confirmed conceptual or constructive (not operational) hypotheses of the form "Construct A is positively related to Construct B in Population C." We are interested in the *construct* "tendency toward eating disorders," not in the *datum* "has reference made to overeating in the nurse's notes for May 15th." An operational replication tests whether we can duplicate our findings using the same methods of measurement and sampling; a constructive replication goes further in the sense of testing the validity of these methods.

Thus, if I cannot confirm Sapolsky's results for patients from my hospital, assessing eating disorders by means of informant interviews, say, or actual measurements of food intake, then clearly Sapolsky has *not* demonstrated any "fact" about eating disorders among psychiatric patients in general. I could then revert to an operational replication, assessing eating problems from the psychiatric notes as Sapolsky did and selecting my sample to conform with the age, sex, and diagnostic properties of his, although I might not regard this endeavor to be worth the effort since, under these circumstances, even a successful operational replication could not establish an empirical conclusion of any great generality or interest. Just as a reliable but invalid test can be said to measure something, but not what it claimed to measure, so an experiment which replicates operationally but not constructively could be said to have demonstrated something, but not the relation between meaningful constructs, generalizable to

some broad reference population, which the author originally claimed to have established.<sup>3</sup>

#### RELATION OF THE SIGNIFICANCE TEST TO THE PROBABILITY OF A "SUCCESSFUL" REPLICATION

The probability values resulting from significance testing can be directly used to measure one's confidence in expecting a "successful" literal replication only. Thus, we can be 98% confident of finding at least 6 of 31 frog responders to have eating problems only if we reproduce all of the conditions of Sapolsky's experiment with absolute fidelity, something that he himself could not undertake to do at this point. Whether we are entitled to anything approaching such high confidence that we could obtain such a result from an operational replication depends entirely upon whether Sapolsky has accurately specified all of the conditions which were in fact determinative of his results. That he did not in this instance is suggested by the fact that, investigating the feasibility of replicating his experiment at the University of Minnesota Hospitals, I found that I should have to review several thousand case records in order to turn up a sample of 31 frog responders like his. Although he does not indicate how many records he examined, one strongly suspects that the base rate of Rorschach frog responding must have been higher at Sapolsky's hospital, either because of some difference in the patient population or, more probably, because an investigator's being interested in some class of responses will tend to subtly elicit such responses at a higher rate unless the testing procedure is very rigorously controlled. If the base rates for frog responding are so different at the two hospitals, it seems doubt-

<sup>3</sup> This distinction between operational and constructive replication seems to have much in common with that made by Sidman (1960) between what he calls "direct" and "systematic" replication. However, in the operant research context to which Sidman directs his attention, "replication" means to run another animal or the same animal again; thus, direct replication involves maintaining the same experimental conditions in detail whereas in systematic replication one allows all supposedly irrelevant factors to vary from one subject to the next in the hope of demonstrating that one has correctly identified the variables which are really in control of the behavior being studied.

ful that the response can have the same correlates or meaning in the two populations and therefore one would be reckless indeed to offer high odds on the outcome of even the most careful operational replication. The likelihood of a successful constructive replication is, of course, still smaller since it depends on the additional assumptions that Sapolsky's samples were truly representative of psychiatric patients in general and that his method of assessing eating problems was truly valid, that is, would correlate highly with a different, equally reasonable appearing method.

#### ANOTHER EXAMPLE

It is not my purpose, of course, to criticize statistical theory or method but rather to suggest ways in which these tools are sometimes misused or misinterpreted by writers or readers of the psychological literature. Nor do I mean to abuse a particular investigator whose research report happened to serve as a convenient illustration of the components of the argument. An abundance of articles can be found in the journals which exemplify these points quite as well as Sapolsky's but space limitations forbid multiple examples. As a compromise, therefore, I offer just one further illustration, showing how the application of these same critical principles might have increased a reader's—and perhaps even an editor's—skepticism concerning some research of my own.

The purpose of the experiment in question (Lykken, 1957) was to test the hypothesis that the "primary" psychopath has reduced ability to condition anxiety or fear. To segregate a subgroup in which such primary psychopaths might be concentrated, I asked prison psychologists to separate inmates already diagnosed as psychopathic personalities into one group that met 14 rather specific clinical criteria specified by Cleckley (1950, pp. 355-392) and to identify another group which clearly did not fit some of these criteria. The normal control subjects were comparable to the psychopathic groups in age, IQ, and sex. Fear conditioning was assessed using the GSR as the dependent variable and a rather painful electric shock as the unconditioned stimulus (UCS). On the index used to measure rate of conditioning, the primary psychopathic

group scored significantly lower than did the controls. By the usual reasoning, therefore, one might conclude that this result demonstrates that primary psychopaths are abnormally slow to condition the GSR, at least with an aversive UCS, and this empirical fact in turn provides significant support for the theory that primary psychopaths have defective fear-learning ability (i.e., a low "anxiety IQ").

But to anyone who has actually participated in research of this kind, this seemingly straightforward reasoning must appear appallingly oversimplified. It is quite impossible to obtain anything resembling a truly random sample of psychopaths (or of nonpsychopathic normals either, for that matter) and it is a matter of unquantifiable conjecture how a sample obtained by a different investigator using equally defensible methods might perform on the tests which I employed. Even with the identical sample, no two investigators are likely to measure the GSR in the same way, use the same conditioned stimulus (CS) and UCS or the same pattern of reinforced and CS-only trials. Given even the same set of protocols, there is no standard formula for obtaining an index of degree or rate of conditioning; the index I used was essentially arbitrary and whether it was a good one is a matter of opinion. My own evaluation of the methods used, together with a complex set of supplementary assumptions difficult to explicate, leads me to believe that these results increase the likelihood that primary psychopaths have slower GSR conditioning with an aversive UCS; I might now give odds of two to one that this empirical generalization is true and odds of three to two that another investigator would be able to confirm it by means of a constructive replication. But this already biased claim is far more modest than the one which is implicit in the significance testing operation, namely, "such a mean difference would only be expected 5 times in 100 if the [generalization] is not true."

This empirical generalization, about GSR conditioning, is derivable from the hypothesis of interest, that psychopaths have a low anxiety IQ, by a chain of reasoning so complex and elliptical and so burdened with accessory assumptions as to be quite impossible to spell

out in the detail required for rigorous logical analysis. Psychologists knowledgeable in the area can evaluate whether it is a reasonable derivation but their opinions will not necessarily agree. Moreover, even if the derivation could pass the scrutiny of some "Certified Public Logician," confirmation of the prediction about GSR conditioning should add only very slightly to our confidence in the hypothesis about fear conditioning. Even if this confirmation were made relatively more firm by, for example, constructive replication of the generalization, "aversive GSR conditioning is retarded in primary psychopaths," the hypothesis that these individuals have a low anxiety IQ could still be said to have passed only the weakest kind of test. This is so because such simple directional predictions about group differences have nearly a 50-50 chance of being true a priori even if our particular hypothesis is false. There are doubtless many possible explanations for low GSR conditioning scores in psychopaths other than the possibility of defective fear conditioning. Indeed, some of my subjects whose conditioning scores were nearly as low as those of the most extreme primary psychopaths seemed to me to be clearly neurotic with considerable anxiety and I attempted to account for their GSR performance with an ad hoc conjecture involving a kind of repression phenomenon, that is, a denial that a low GSR index implied poor fear conditioning in their cases.

A redeeming feature of this study was that two other related but distinguishable predictions from the same hypothesis were tested at the same time, namely, that primary psychopaths should do as well as normals on a learning task involving positive reward but less well on an avoidance learning problem, and that they should be more willing than normals to choose embarrassing or frightening situations in preference to alternatives involving tedium, frustration, physical discomfort, and the like. Tests of these predictions gave affirmative results also, thus providing some of the multiple corroboration necessary for the hypothesis to claim the attention of other experimenters.

Obviously, I do not mean to criticize the editor's decision to publish my (1957) paper. The tendency to evaluate research in terms of

mechanical rules based on the results of the significance tests should not be replaced by equally rigid requirements concerning replication or corroboration. This study, like Sapolsky's or most others in this field, can be properly evaluated only by a qualified reader who can substitute his own informed judgment and scientific intuition for the rigorous reasoning and experimental control that is usually not achievable in clinical and personality research. As it happens, subsequent work has provided some encouraging support for my 1957 findings. The two additional predictions mentioned above have received operational replication (i.e., the same test methods used in a different context) by Schachter and Latené (1964). The prediction that psychopaths show slower GSR conditioning with an aversive UCS has been constructively replicated (i.e., independently tested with no attempt to copy my procedures) by Hare (1965a). Finally, two additional predictions from the theory that the primary psychopath has a low anxiety IQ have been tested with affirmative results (Hare, 1965b; 1966). All told, then, this hypothesis can now boast of having led to at least five quasi-independent predictions which have been experimentally confirmed and three of which have been replicated. The hypothesis is therefore entitled to serious consideration although one would be rash still to regard it as proven. At least one alternative hypothesis, that the psychopath has an unusually efficient mechanism for inhibiting emotional arousal, can account equally well for the existing findings so that, as is usually the case, further research is called for.

#### CONCLUSIONS

The moral of this story is that the finding of statistical significance is perhaps the least important attribute of a good experiment; it is *never* a sufficient condition for concluding that a theory has been corroborated, that a useful empirical fact has been established with reasonable confidence—or that an experimental report ought to be published. The value of any research can be determined, not from the statistical results, but only by skilled, subjective evaluation of the coherence and reasonableness of the theory,

the degree of experimental control employed, the sophistication of the measuring techniques, the scientific or practical importance of the phenomena studied, and so on. Ideally, all experiments would be replicated before publication but this goal is impractical. "Good" experiments will tend to replicate better than poor ones (and, when they do not, the failures will tend to be informative in themselves, which is not true for poor experiments) and should be published so that they may stimulate replication and extension by others. Editors must be bold enough to take responsibility for deciding which studies are good and which are not, without resorting to letting the  $p$  value of the significance tests determine this decision. There is little real danger that anything of value will be lost through this approach since the unpublished investigator can always resort to constructive replication to induce editorial acceptance of his empirical conclusions or to multiple corroboration to compel editorial respect for his theory. Since operational replication must really be done by an independent second investigator and since constructive replication has greater generality, its success strongly implying that an operational replication would have succeeded also, one should usually replicate one's own work constructively, using different sampling and measurement procedures within the purview of the same constructive hypothesis. If only unusually well done, provocative, and important research were published without such prior authentication, operational replication of such research by others

would be come correspondingly more valuable and entitled to the respect now accorded capable replication in the other experimental sciences.

## REFERENCES

- BAKEN, D. The test of significance in psychological research. *Psychological Bulletin*, 1966, **66**, 423-437.
- CLECKLEY, H. *The mask of sanity*. Saint Louis: C. V. Mosby, 1950.
- HARE, R. D. Acquisition and generalization of a conditioned fear response in psychopathic and non-psychopathic criminals. *Journal of Psychology*, 1965, **59**, 367-370. (a)
- HARE, R. D. Temporal gradient of fear arousal in psychopaths. *Journal of Abnormal Psychology*, 1965, **70**, 442-445. (b)
- HARE, R. D. Psychopathy and choice of immediate versus delayed punishment. *Journal of Abnormal Psychology*, 1966, **71**, 25-29.
- HOLTZMAN, W. H., & BITTERMAN, M. E. A factorial study of adjustment to stress. *Journal of Abnormal and Social Psychology*, 1956, **52**, 179-185.
- LYKKEN, D. T. A study of anxiety in the sociopathic personality. *Journal of Abnormal and Social Psychology*, 1957, **55**, 6-10.
- MEEHL, P. E. Theory-testing in psychology and physics: A methodological paradox. *Philosophy of Science*, 1967, **34**, 103-115.
- ROSENTHAL, R. On the social psychology of the psychological experiment: The experimenter's hypothesis as unintended determinant of experimental results. *American Scientist*, 1963, **51**, 268-283.
- SAPOLSKY, A. An effort at studying Rorschach content symbolism: The frog response. *Journal of Consulting Psychology*, 1964, **28**, 469-472.
- SCHACHTER, S., LATENÉ, B. Crime, cognition and the autonomic nervous system. *Nebraska Symposium on motivation*, 1964, **12**, 221-273.
- SIDMAN, M. *Tactics of scientific research*. New York: Basic Books, 1960.

(Received March 8, 1967)