# ON THE REJECTION OF OUTLIERS

THOMAS S. FERGUSON

UNIVERSITY OF CALIFORNIA, LOS ANGELES AND PRINCETON UNIVERSITY

## 1. Introduction and summary

The general problem treated in this paper is a very old and common one. In its simplest form it may be stated as follows. In a sample of moderate size taken from a certain population, it appears that one or two values are surprisingly far away from the main group. The experimenter is tempted to throw away the apparently erroneous values, and not because he is certain that the values are spurious. On the contrary, he will undoubtedly admit that even if the population has a normal distribution there is a positive although extremely small probability that such values will occur in an experiment. It is rather because he feels that other explanations are more plausible, and that the loss in the accuracy of the experiment caused by throwing away a couple of good values is small compared to the loss caused by keeping even one bad value. The problem, then, is to introduce some degree of objectivity into the rejection of the outlying observations.

There is no need to give here a historical outline of progress in the subject. Good accounts of the historical aspect, interesting because this subject was one of the first problems to receive a statistical treatment, may be found in recent papers by Grubbs [4], Murphy [9], and Anscombe [1]. We shall mention here only those papers which directly concern us, with particular emphasis placed on the last ten years.

Two mathematical models have been proposed, implicitly by Grubbs and explicitly by Dixon [3], to give a structure to the outlier problem. In both models it is assumed that a sample of n observations, to be denoted by $X_1, \cdots, X_n$ has been drawn from a normal population with possibly unknown mean and variance. A few of these values may have been spuriously changed. In *model A*, this change is hypothesized as a shift in the mean, and in *model B*, as an increase in the variance. The precise formulations will be made clear in the particular problems we shall consider.

This paper is mainly concerned with the derivation, found in section 2, of the locally best tests for the existence of spurious observations in several situations for both models A and B. The tests suggested here by virtue of their local optimality are based on the sample coefficient of skewness for one-sided alternatives, and on the sample coefficient of kurtosis for two-sided alternatives. Many spurious observations are allowed under the alternative hypotheses; even in the most