



**HAL**  
open science

# OUTLIER DETECTION IN GRIDDED SHIP'S DATASETS

Pascal Terray

► **To cite this version:**

Pascal Terray. OUTLIER DETECTION IN GRIDDED SHIP'S DATASETS. Advances in the applications of marine climatology- The Dynamic Part of the WMO Guide to the Applications of Marine Climatology, WMO/TD-No. 1081 (13), pp.177-186, 2003, JOINT WMO/IOC TECHNICAL COMMISSION FOR OCEANOGRAPHY AND MARINE METEOROLOGY TECHNICAL REPORT SERIES. hal-01413807

**HAL Id: hal-01413807**

**<https://hal.science/hal-01413807v1>**

Submitted on 11 Dec 2016

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# OUTLIER DETECTION IN GRIDDED SHIP'S DATASETS

Pascal Terray

*Laboratoire d'Océanographie Dynamique et de Climatologie, and Université Paris 7, Paris, France.*

## 1. INTRODUCTION

This is the second of two papers attempting to develop robust statistical methods to deal with gridded ship's datasets. The earlier study (Terray, 1998) focused on an extension of the traditional Empirical Orthogonal Function (EOF) analysis which allows arbitrary positive weights to be assigned to each entry of the data matrix. If these weights are constructed in a responsible manner (for example, as a smooth function of the number of ship reports used to compute a particular raw monthly mean in the dataset), it was demonstrated that this method allows us to analyze the natural variability exhibited by gridded ship's datasets by directly taking into account the irregular space-time sampling of marine observations. In particular, the method takes care of missing values by assigning zero weights to such data entries.

In the current study, we discuss another robust statistical method to detect « local errors » in gridded ship's datasets. More precisely, we attack the problem of outlying areal averages in gridded ship's datasets such as Comprehensive Ocean-Atmosphere Data Set (COADS; Woodruff et al. 1987)  $2^\circ$  lat  $\times$   $2^\circ$  long monthly summaries and how to test the statistical significance of them. Since the majority of climate researchers use gridded ship's datasets instead of individual ship reports, we suggest that these datasets must be checked for the presence of doubtful raw monthly means in the same manner as individual ship reports are quality controlled before being integrated in ship reports databases. Moreover, it should be noted that such approach may be a solution to the trimming problems which are apparent in COADS monthly summaries (Wolter, 1997).

The rest of this paper is organized as follows: first, we present some elements of outlier detection theory and the basic statistical tests we have used. Next, we discuss how these statistical tests may be adapted to ship's datasets and integrated as building blocks in a fully computerized procedure for detecting many outliers in such datasets. Finally, this new approach has been experimented on a marine product in order to show how it works in practice. As a conclusion, we suggest that the two procedures, outlier detection and weighted EOF analysis, may be combined to obtain a truly robust statistical method particularly well suited to gridded ship's datasets.

## 2. STATISTICAL THEORY OF OUTLIER DETECTION

In the context of gridded ship's datasets, an outlying observation, or "outlier", is a raw monthly mean in a  $2^\circ$  lat  $\times$   $2^\circ$  long box (depending on the resolution of the dataset) that appears to deviate markedly from adjacent or neighboring grid-points in area or/and in time. Outliers in gridded ship's datasets may be generated by three basic mechanisms (Wolter, 1997):

- An outlying raw monthly mean may be merely an extreme manifestation of the sampling inherent in the data, since some raw monthly means in  $2^\circ$  lat  $\times$   $2^\circ$  long boxes are computed with very few marine observations for a given date while adjacent boxes may be well sampled.
- outlying raw monthly means in some  $2^\circ$  lat  $\times$   $2^\circ$  long boxes may also be the results of potential biases due to the origin of the "source-decks" merged into the gridded ship's dataset or processing errors. For example, biases in Sea Surface Temperature (SST) associated with different methods of measurements (bucket or intake) may well introduce errors in gridded ship's datasets in particular atmospheric conditions and along some ship's tracks.
- Finally, an outlying areal average may be the result of errors relating to instrumental readings or coding mistakes. But, most of these types of outliers must be discovered during basic quality controls which are automatically applied to individual ship reports merged into any reasonable marine product.

The problem of detecting outliers in random sample has been extensively researched by statisticians in recent years and a number of test statistics are available for both the single outlier case and the many outlier case for testing a specified number  $k$  of outliers (Barnett and Lewis, 1978). In particular, the detection of outliers in normal sample has received considerable attention. It is far beyond the scope of this paper to give a review of this vast subject. Suffice to say here, that the problem of outlier detection is generally treated as the statistical testing of a hypothesis. The null hypothesis, as usually stated, is that all the observations are drawn from the same (normal) population; the alternative hypothesis is that at least one of the observations has been drawn from an other distribution. In order to discriminate between these two hypotheses, a sample criterion  $T$  which uses the doubtful

observation(s) is calculated. This statistic is then compared with a critical value  $\lambda_{\alpha}$  based on the theory of random sampling to determine whether the doubtful observation is to be retained or rejected. This critical value is the value of the chosen sample criterion which would be exceeded by chance with some specified and small probability  $\alpha$  (say 0.01 or 0.05), the so-called significance level of the test, if the null hypothesis is true. Intuitively, this significance level is the risk or erroneously rejecting a good observation (statistical type I error). More precisely, statistical tests for outliers are then of the following form:

- 1) Find  $\lambda_{\alpha}$  such that  $\Pr(T > \lambda_{\alpha}) = \alpha$  if the null hypothesis is true for some statistic T;
- 2) Reject the null hypothesis and declare an outlier present if  $T > \lambda_{\alpha}$  or accept the null hypothesis and declare the sample is clean if  $T \leq \lambda_{\alpha}$ .

In this statistical framework, outlier detection procedures differ by:

- the form of the underlying parent population (normal, gamma, ... );
- the form of the test criterion T which has to be computed on the sample: among these test criteria, we can distinguish those which clearly identify particular observations as possible outliers from those which test the hypothesis that the random sample as a whole did indeed come from the specified parent distribution;
- the number of suspected outliers in the sample;
- the fact that the doubtful observations may be to one side of the bulk of the data or that some are too large and some are too small.

Several hundreds of statistical tests of this type are described in the book of Barnett and Lewis (1978) which is a kind of "bible" on the subject. In the context of gridded ship's datasets, the problem is then to decide which tests to apply and how to use them in order to obtain a fully computerized procedure for detecting outliers which may be applied to any ship's dataset. In this way, one can hope to trap anomalous cases and so ensure the integrity of most of the ship's datasets currently in use.

We have used here a simple model, well documented in the statistical literature: when the data with the possible exception of any outlier, form a sample from a normal distribution with unknown mean  $\mu$  and unknown variance  $\sigma^2$ . We recognize that this model is certainly not perfect in the context of samples of adjacent raw monthly means in  $2^{\circ}$  lat  $\times$   $2^{\circ}$  long boxes extracted from gridded ship's datasets. However, as we will show below, this model works "reasonably" well as implemented in our computerized procedure on the basis of the spatial coherence of neighboring  $2^{\circ}$  lat  $\times$   $2^{\circ}$  long area values for many meteorological parameters.

Several statistical tests exist that are reasonably powerful to detect one outlier in a normal sample and our approach involves the following classical statistical criteria:

Let  $x_1, x_2, \dots, x_n$  be the observations of a random sample. Order the observations according to increasing magnitude and denote the  $i$ th largest by  $y_i$ ; thus,  $y_1 \leq y_2 \leq \dots \leq y_n$  is the ordered set of observations. Suppose the largest observation  $y_n$  is suspect. In order to test for discordancy this single upper observation in a normal sample, a reasonable test statistic is

$$T = \frac{y_n - \bar{x}}{s}$$

where

$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n}$$
 is the sample mean,

$$s^2 = \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{n} \right)^2$$
 is the sample variance

calculated with  $n$  degrees of freedom.

If  $y_1$ , the lower observation, rather than  $y_n$  is the doubtful value, the criterion is as follows:

$$T' = \frac{\bar{x} - y_1}{s}$$

and the rest of the statistical procedure will be unchanged on the basis of the symmetry of the normal distribution. Finally, when it is not known a priori whether the contaminant is the lower or the upper observation in the sample, we should compute

$$T^* = \max(T, T')$$

But, in this last case, we must use a critical value corresponding to the  $\alpha/2$  significance level if we want the true significance level to be 0.05.

The rationale behind these tests may be found in Hawkins (1980) or Barnett and Lewis (1978). The null hypothesis that we are testing in every cases is that all the observations in the sample come from the same normal population. It may be shown that these statistics are optimal in the sense of maximizing the probability of correct identification of an outlier when one is present. It should be noted, however, that these statistics may produce quite misleading results in the presence of many outliers, especially when suspected values are closer to each other than they are close to the bulk of the other observations. This inability of a testing procedure to identify even a single outlier in the presence of several suspected values is called the masking effect. We will discuss further this point in the next section

when we describe our computerized procedure for detecting outliers.

Before using these test statistics in outlier checks, we must know the significance probability attaching to an observed value  $t$  of the statistic  $T$  (or  $T^*$ ). That is to say, the probability that, on the null hypothesis of no contamination,  $T$  takes values more discordant than  $t$ . For this purpose, we need to find the null distribution of  $T$  or at least some fractiles  $\lambda_\alpha$  of this distribution corresponding to specified significance levels  $\alpha$ , say 0.01 or 0.05. The null distribution of  $T$  is available as a recursion relationship (Barnett and Lewis, 1978) or as a complicated multiple integral (Grubbs, 1950), and tables containing critical values for some standard significance levels have been published (Grubbs and Beck, 1972; Hawkins, 1980). However, we will show how approximate critical values for a given significance level  $\alpha$  can be computed since our computerized procedure may involve number of observations outside of the range of these published tables.

Without loss of generality, we consider only the case of an upper outlier, approximate fractiles for  $T$  or  $T^*$  may be derived similarly. We may compute some fractiles of the test distribution of  $T$  as follows:

Under the null hypothesis of no contamination,  $x_1, x_2, \dots, x_n$  are observations of random variables  $X_1, X_2, \dots, X_n$  which are independent and identically distributed as  $N(\mu, \sigma^2)$ . In this case, if  $x_i$  is an observation selected arbitrarily from the random sample of  $n$  items, it may be shown that if

$$T_i = \frac{x_i - \bar{x}}{s}$$

then the probability density function of

$$t_i = \frac{T_i \sqrt{n-2}}{\sqrt{n-1-T_i^2}}$$

is given by the "Student's"  $t$ -distribution with  $n-2$  degrees of freedom. This is easily verified because  $t_i$  is the test statistic of the classical Student's two sample  $t$ -test, where one sample consists of  $x_i$  and the second sample of the  $n-1$  other observations. From this result, we are able to find the probability that an arbitrary observation  $i$  will be outlying since

$$\Pr[T_i > \lambda] = \Pr\left[t_{(n-2)} > \frac{\lambda \sqrt{n-2}}{\sqrt{n-1-\lambda^2}}\right]$$

where  $\lambda$  is an arbitrary value in the range  $[-\sqrt{n-1}, \sqrt{n-1}]$  and  $t_{(n-2)}$  follows a Student's  $t$ -distribution with  $n-2$  degrees of freedom. However, this result does not yet give us an exact test for one outlier, because this probability is different from the probability that a particular observation (the lowest or the largest) will be greater than  $\lambda$ . More precisely, we need the distribution not of an arbitrary  $T_i$ , but of  $T$ , the greatest of the quantities  $T_i$  for  $i=1$  to  $n$ .

Now, note that the event  $(T > \lambda)$  is the union of the  $n$  events  $(T_i > \lambda)$ . Thus,

$$\Pr[T > \lambda] = \Pr\left[\bigcup_{i=1}^n (T_i > \lambda)\right]$$

In words, the probability of the event  $(T > \lambda)$  is the probability that at least one of the  $n$  events  $(T_i > \lambda)$  is true. Bounds on  $\Pr[T > \lambda]$  may then be obtained in terms of the component events  $(T_i > \lambda)$  through the use of the so-called Bonferroni inequality (Feller, 1968)

$$\begin{aligned} \sum_i \Pr[T_i > \lambda] - \sum_{i < j} \Pr[(T_i > \lambda) \cap (T_j > \lambda)] &\leq \\ \Pr[T > \lambda] &\leq \sum_i \Pr[T_i > \lambda] \end{aligned}$$

Since the events  $(T_i > \lambda)$  are equiprobable, and likewise the events  $(T_i > \lambda) \cap (T_j > \lambda)$ , we have the following inequality for arbitrary  $i$  and  $j$

$$\begin{aligned} n \Pr[T_i > \lambda] - \frac{n(n-1)}{2} \Pr[(T_i > \lambda) \cap (T_j > \lambda)] &\leq \\ \Pr[T > \lambda] &\leq n \Pr[T_i > \lambda] \end{aligned}$$

Now, by using the fact that for arbitrary  $i$  and  $j$  (Doornbos, 1966)

$$\Pr[(T_i > \lambda) \cap (T_j > \lambda)] < (\Pr[T_i > \lambda])^2$$

we finally obtain

$$\begin{aligned} n \Pr[T_i > \lambda] - \frac{n-1}{2n} (n \Pr[T_i > \lambda])^2 &< \Pr[T > \lambda] \\ &\leq n \Pr[T_i > \lambda] \end{aligned}$$

for an arbitrary  $i$ . Thus, if

$$\frac{\lambda \sqrt{n-2}}{\sqrt{n-1-\lambda^2}}$$

is the  $1 - (\alpha/n)$  fractile of the Student's t-distribution with  $n-2$  degrees of freedom, the last equation shows that

$$\alpha - \frac{n-1}{2n}(\alpha^2) < \Pr[T > \lambda] \leq \alpha$$

A result indicating that  $\lambda$  is a good and conservative approximation of the true critical value  $\lambda_\alpha$  of the distribution of  $T$  under the null hypothesis of no contamination for any reasonable significance level  $\alpha$ , say 0.01 or 0.05. Moreover, it can be shown that this method gives the exact critical value  $\lambda_\alpha$  of  $T$  if

$$\lambda \geq \sqrt{\frac{n-1}{2}}$$

(for example, the 0.05 critical value for any  $n < 15$ ), since in this case we have

$$\Pr[(T_i > \lambda) \cap (T_j > \lambda)] = 0$$

for arbitrary  $i$  and  $j$ . Following the same procedure, we may approximate the true critical value  $\lambda_\alpha$  of  $T^*$  on the null hypothesis of no contamination by  $\lambda^*$ , if

$$\frac{\lambda^* \sqrt{n-2}}{\sqrt{n-1 - (\lambda^*)^2}}$$

is the  $1 - \alpha/(2n)$  fractile of the Student's t-distribution with  $n-2$  degrees of freedom.

### 3. OUTLIER DETECTION IN GRIDDED SHIP'S DATASETS

Suppose now that we want to check the « local » consistency of a given ship's dataset with, say, a  $2^\circ$  lat  $\times$   $2^\circ$  long resolution. In our approach, the preceding theoretical results are then used as follows:

- 1) First, we specify upper and lower limits for detecting doubtful raw monthly mean values in the gridded ship's dataset. These limits may vary depending on calendar month and area. Any value which exceeds the upper limit or is less than the lower limit is considered a priori doubtful and will be tested for compatibility with raw monthly mean values in adjacent or neighboring  $2^\circ$  lat  $\times$   $2^\circ$  long boxes.
- 2) For any date, doubtful raw monthly mean values identified in step 1 are arranged from the most outlying to the most inlying compared to the bulk of the data. For this purpose, absolute values of residuals of these doubtful values from the overall mean of the observed data for this date

are sorted in descending order and the doubtful values are ranked accordingly.

- 3) These doubtful values are then considered consecutively, from the most outlying to the most inlying, and a sample is constructed from adjacent or neighboring  $2^\circ$  lat  $\times$   $2^\circ$  long area values for any of these possible outliers. The number of  $2^\circ$  lat  $\times$   $2^\circ$  long boxes in the vicinity of each doubtful value which are scanned, in order to construct a sample, may be chosen by the user before running the procedure. It should be noted that the number of items in this sample may vary depending on the date and the area. However, the significance level  $\alpha$  of the test will be the same for any suspected raw monthly mean value, as we will see below.

- 4) At this stage, several different possibilities exist :

- a) First, we need to consider the case when it is not possible to pick up a sample to test the doubtful value because none surrounding boxes contain data. Frequently, this means that the doubtful raw monthly mean is calculated from very few ship reports. In such case, the user may decide, before running the procedure, to flag or to reject all these unrepresentative values.

- b) Second, suppose that there is only one doubtful value in the collected sample, the one we want to test. If this value is at the upper end of the sample, we use  $T$  as a test criterion; if it is at the lower end of the sample, the statistic  $T'$  is considered instead. In both cases, the doubtful value is declared as an outlier if the statistic exceeds the critical value  $\lambda_\alpha$  corresponding to a specified significance level  $\alpha$ . In this case, the suspected value is rejected or flagged (a user choice) and the next most outlying doubtful value is processed.

- c) Finally, imagine there is more than one doubtful value in the constructed sample of  $n$  items, according to the upper and lower limits specified in step 1. Let  $K$  be the number of such doubtful values and  $x$  the suspected value we are currently processing. In order to take into account the possibility that the sample contains more than one outlier, a consecutive procedure is applied. One naive approach is to use repeated applications of the single outlier statistical test  $T^*$  described above, deleting the "outlier" detected at each step and applying the test again to the reduced sample until an insignificant result is obtained or the suspected value  $x$  is tested for compatibility with the remaining observations. However, this "forward selection" approach may be quite misleading in practice (Hawkins, 1980). The problem is the so-called masking effect discussed in the preceding section: the presence of two or more outliers may produce an insignificant result in the initial single outlier

test. In view of this defect, the following variant is recommended : remove the  $K$  most extreme values of the sample (absolute values of residuals from the sample mean of the successively reduced sample are used to rank the observations). If the current doubtful value  $x$  is not thrown away in this process, declare  $x$  has « clean » and process the next most outlying doubtful value for the current date. Otherwise, apply the following "backward selection" algorithm: starting with the  $n-K$  « clean » observations, test the most inlying of the  $K$  extreme values for compatibility with the clean observations by the statistic  $T^*$  at a nominal significance level  $\alpha$ . If it is compatible, then include it with the clean observations and repeat the procedure with the next most outlying suspected value in the sample until the current doubtful value  $x$  is processed and is declared as compatible. This sequence of tests is immediately stopped when an observation is rejected by the statistical test  $T^*$  since all the subsequent outlying raw monthly means, including  $x$ , are then incompatible with the clean observations. In this case, the  $2^\circ$  lat  $\times$   $2^\circ$  long area mean value corresponding to  $x$  is rejected or flagged, and the next doubtful value for the current date is processed. Note, however, that the other rejected values in the sample are not set to missing at this stage.

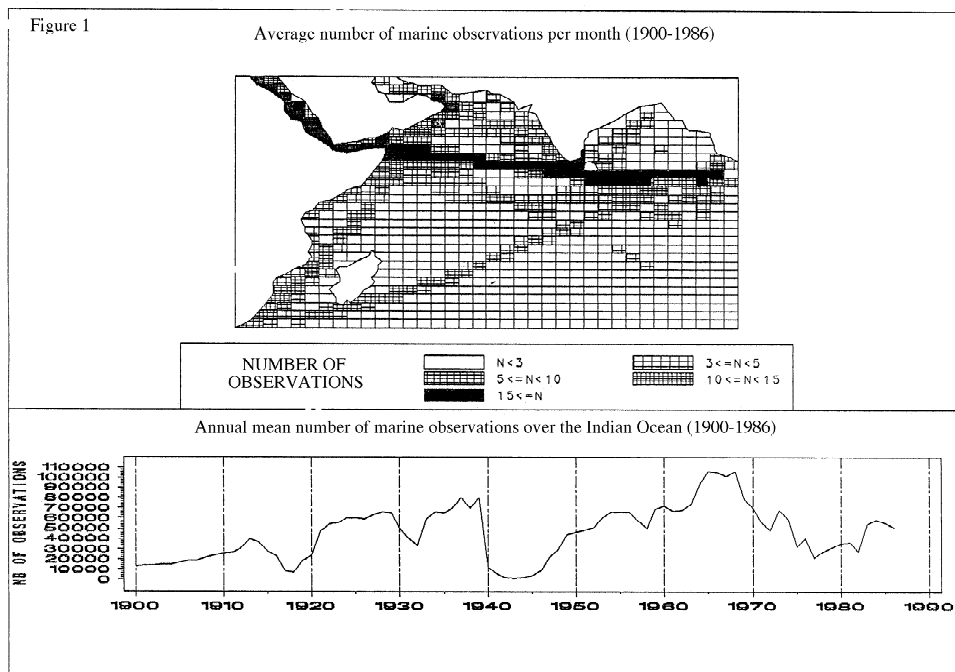
It is fair to say that, while the backward consecutive algorithm described in 4.c is immune to masking (providing that the actual number of outliers in the sample does not exceed the number of suspected values  $K$  in the test procedure), it provides important distributional difficulties associated with finding suitable fractiles  $\lambda_\alpha$  if we require (as we do) an actual significance level  $\alpha$  for each of the

successive null hypotheses which are tested in the backward selection algorithm. A comprehensive discussion of this problem is given by Hawkins (1980) and we omit the details due to the lack of space and difficulty of the problem. Suffice to say here, that it is necessary to resort to simulation if we require exact fractiles, but that there is little error by approximating these fractiles as outlined in the preceding section excepted for small  $n$ , say  $n < 15$ . The latter solution was adopted in this study. The consequence of this is that the sequence of tests used in the backward consecutive algorithm described in 4.c may have actual significance levels in excess of 25% of the specified nominal significance level  $\alpha$  according to Hawkins (1980). We will try to correct this deficiency in a future version of our outlier detection procedure by doing the required simulations.

#### 4. EXAMPLE

The outlier detection algorithm has been applied to several ship's datasets and various examples will be presented during the workshop. In particular, an experiment was undertaken on a pre-COADS marine product with known systematic errors, in order to show the benefit of this type of procedure in the context of marine climatology.

An extensive description of the ship's dataset used in this experiment may be found in Terray (1994). Briefly, SST data are presented as raw monthly means in  $2^\circ$  lat  $\times$   $2^\circ$  long boxes in a domain extending from  $30^\circ$  to  $100^\circ$ E Longitude and from  $30^\circ$ S to  $30^\circ$ N Latitude. The period of analysis extends from 1900 to 1986. Figure 1 documents the irregular space-time sampling associated with this gridded ship's dataset.



Many well-known deficiencies were observed in this dataset before and around World War II (Terray, 1994). In addition, a suspicious warming trend is apparent on the SST time series during 1954-1976 and it was anticipated that this trend may be linked to important changes in the origin of “source-decks” merged into this marine product or to the presence of an huge amount of erroneous ship-reports which were not rejected during basic quality control of the ship reports. Suspect raw monthly means were mainly confined along the shipping routes from Madagascar to Sumatra and from Sumatra to the Northern Arabian Sea for the period 1968-1974.

In view of this, the outlier detection algorithm of the preceding section has been applied to this SST gridded ship’s dataset in a two steps procedure :

- First, the algorithm was applied to all the raw monthly SST fields with 15°C as a lower limit and 35°C as an upper limit to identify doubtful 2° lat × 2° long monthly means which must be tested by the algorithm. A nominal significance level of 0.05 was chosen for all the tests. This first step was only intended as a check on “evident” outliers far away from the bulk of the data. 481 raw monthly values were tested in this first step for all the monthly fields of 1900-1986 and, among them, 361 were identified as outliers by the statistical tests (this number includes isolated monthly mean values) and were rejected.
- The second step is designed to remove outliers with respect to anomaly fields. For this purpose, the raw monthly mean SST fields were expressed as monthly anomaly fields by using a monthly climatology obtained from a weighted EOF analysis on COADS SST data (Terray, 1998). The outlier detection algorithm was applied to these anomaly fields with -3°C as a lower limit and 3°C as an upper limit. Again, a nominal significance level of 0.05 was chosen for all the tests. 10917 anomaly values were tested in this second step; among them, 2826 were identified as outliers and the corresponding raw monthly means values were rejected.

In order to investigate the impacts of the outlier detection algorithm, the following computations were undertaken on the SST ship’s dataset both before and after the “cleaning” of the data:

(i) First, the 1954-1976 interval was used as a reference period for calculating a climatology for each calendar month and each 2° box, provided that data for at least 10 years with more than 5 observations per month were available in the period. The monthly means for each i grid point and j month were computed as a weighted average

$$\bar{X}_{ij} = \frac{\sum_{k=1954}^{1976} W_{ijk} X_{ijk}}{\sum_{k=1954}^{1976} W_{ijk}}$$

where  $W_{ijk} = 1 - \exp(-N_{ijk} / 5)$

Here  $X_{ijk}$  is the value computed for the ith box, jth month and kth year.  $N_{ijk}$  is the number of ship-observations used in computing  $X_{ijk}$ .  $W_{ijk}$  is in the neighborhood of 1 if  $N_{ijk} > 10$  and near 0.5 if  $N_{ijk}$  equals 5.

(ii) After this first step, time monthly anomaly series for each 2° box during the 1900-1986 period were computed by simply subtracting from each value this climatology, provided that neither the datum nor the climatology was missing. These anomalies were then subsequently spatially averaged over the whole Indian Ocean with the same weighting scheme (e.g.,  $W_{ijk}$ ) as used in the computation of the climatology.

The two SST anomaly time series computed, respectively, before and after the « cleaning » of the data, were then subjected to the X11 monthly additive scheme (Terray, 1994), a powerful technique for describing a time series, in order to assess their consistency. In the X11 procedure, the analyzed  $X_t$  monthly time series is decomposed into three terms

$$X_t = T_t + A_t + I_t$$

The  $T_t$  term is used to quantify the trend and low-frequency variations in the time series. The  $A_t$  term describes the annual cycle and the  $I_t$  can be used to assess the level of noise in the data, though this term can also contain some signal in a climatological sense. All the terms are estimated with specific moving averages of various lengths. Figures 2 and 3 give the results of the analysis for the SST time series computed before and after outliers were rejected, respectively. The monthly number of observations is also plotted on the bottom of each figure as an aid for interpreting the results and detecting accurately any change in the composition of the "source-decks" contributing to the time series. While the two series and their associated X11 components are similar in many aspects, an important discrepancy may be noted during 1968-1974: the unlikely warm anomalies observed on the data before running the outlier detection procedure (Figure 2) are considerably reduced on the time series computed after outliers were rejected (Figure 3). As a consequence, the trend components of the two series are different during 1968-1974. This difference is consistent with the hypothesis of the artificial nature of the warming trend observed during 1968-1974 over the Indian Ocean. Finally, it may be noted that the « clean » series is less noisy as demonstrated by the irregular components.

Figure 2: SST monthly anomalies relative to 1954-1976 for the whole Indian Ocean before outlier detection. The series has been decomposed into annual, trend and irregular components by the X11 procedure. The monthly number of ship reports used to construct the series is given at the bottom of the figure.

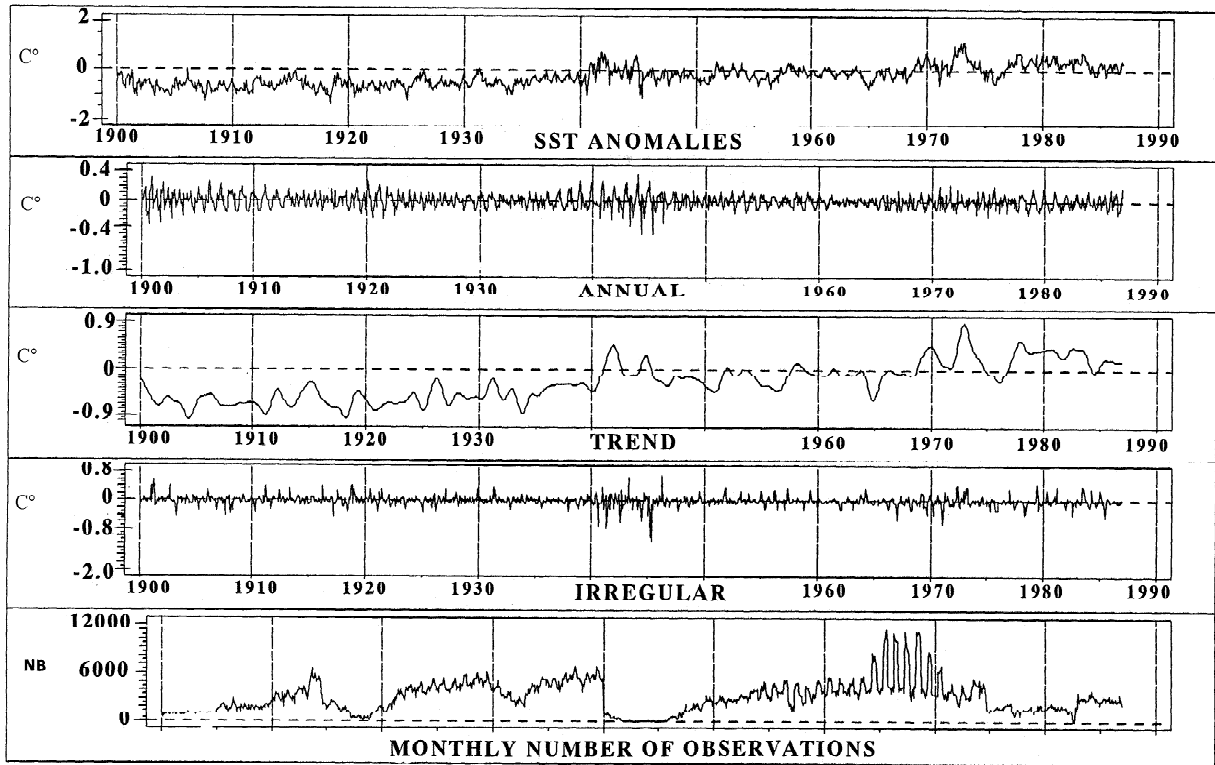
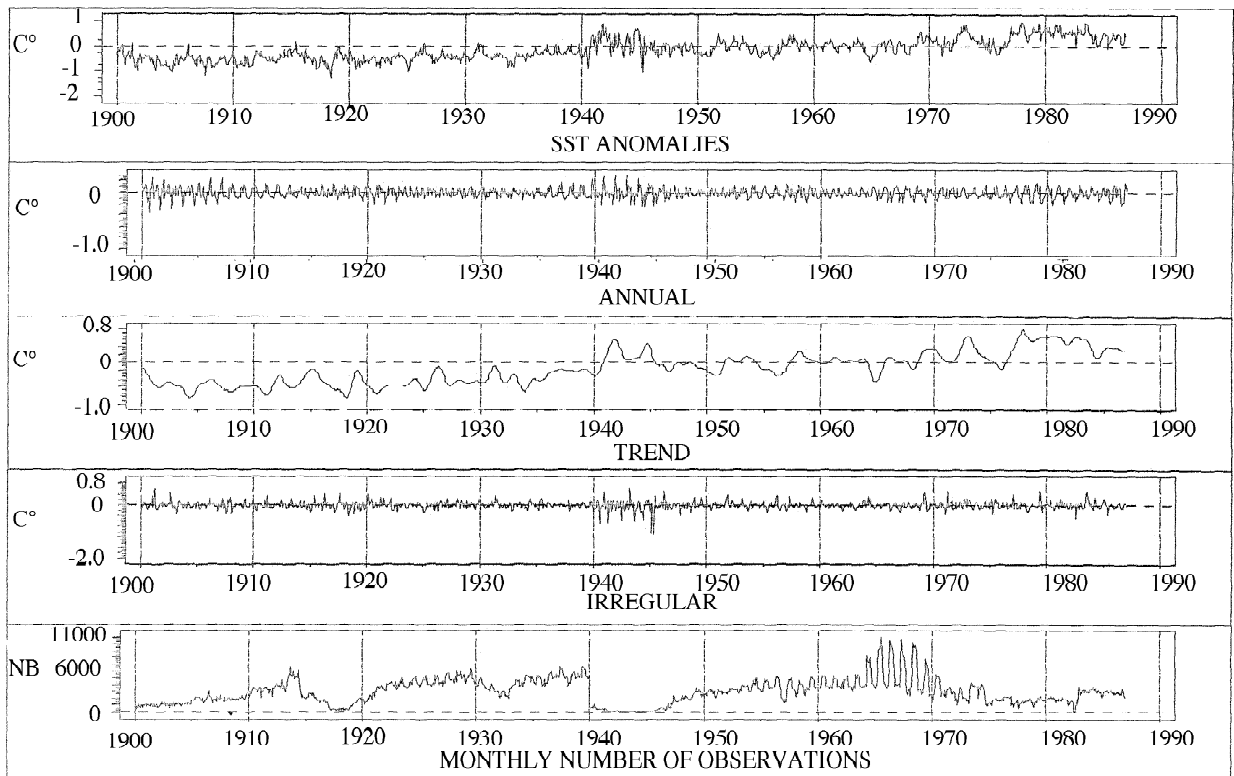


Figure 3: SST anomalies relative to 1954-1976 for the whole Indian Ocean after outlier detection. The series has been decomposed into annual, trend and irregular components by the X11 procedure. The monthly number of ship reports used to construct the series is given at the bottom of the figure.





## 5. CONCLUSIONS

A recurring problem in the creation and maintenance of large gridded ship's datasets is the accuracy of the information entering these products. The fact that large volumes of data are involved suggests that, as far as possible, the reliability of such datasets should be assessed through a computerized screening procedure. For this purpose, a new method for detecting outliers in gridded ship's datasets has been proposed. It is our hope that this approach will aid climate scientists in determining which, if any, of the raw monthly area values included in a particular ship's dataset may be outliers.

Once potential outliers have been identified, it is suggested that these values may be flagged or, more drastically, rejected. In any case, the impact of these doubtful values in a particular data analysis may be easily assessed by comparing the results obtained before and after these potential outliers are rejected. In this way, it may be possible to obtain more reliable results in marine climatology.

The proposed approach may also be considered as a valuable alternative to trimming procedures which are applied to ship reports before computing monthly mean summaries for  $2^\circ$  lat  $\times$   $2^\circ$  long boxes in order to reduce erroneous data losses.

## REFERENCES

- Barnett, V., and T. Lewis, 1978: *Outliers in Statistical Data*. John Wiley & Sons, Inc., New York.
- Doornbos, R., 1966: *Slippage Tests*. 1st edn. Mathematical Centre Tracts, no 15, Mathematisch Centrum, Amsterdam.
- Feller, W., 1968: *An Introduction to Probability Theory and Its Applications*. vol. 1, 3d ed., John Wiley & Sons, Inc., New York.
- Grubbs, F. E., 1950 : Sample criteria for testing outlying observations. *Ann. Math. Statist.*, 21, 27-58.
- Grubbs, F. E., and G. Beck, 1972: Extension of sample sizes and percentage points for significance tests of outlying observations. *Technometrics*, 14, 847-854.
- Hawkins, D.M., 1980: *Identification of Outliers*. Chapman and Hall, London.
- Terray, P., 1994: An evaluation of climatological data in the Indian Ocean area. *J. Meteor. Soc. Japan*. 72, 359-386.
- Terray, P., 1998: Detecting Climatic Signals from Ship's Datasets. Proceedings of International Workshop on Digitization and Preparation of Historical Surface Marine Data and Metadata, 15-17 September 1997, Toledo, Spain. H.F. Diaz and S.D. Woodruff Eds., WMO publication (in press).
- Wolter, K., 1997: Trimming problems and remedies in COADS. *J. Climate*, 10, 1980-1997.
- Woodruff, S.D., R.J. Slutz, R.L. Jenne, and P.M. Steurer, 1987: A Comprehensive Ocean-Atmosphere Data Set. *Bull. Amer. Meteor. Soc.*, 68, 1239-1250.