

## Systematic error of dimension estimates using fixed mass scaling methods

Machiel de Rover<sup>1</sup> and Willem van de Water<sup>2</sup>

<sup>1</sup>*FOM—Instituut voor Plasmafysica “Rijnhuizen,” P. O. Box 1207, 3430 BE Nieuwegein, The Netherlands*

<sup>2</sup>*Physics Department, Eindhoven University of Technology, P.O. Box 513, 5600 MB, Eindhoven, The Netherlands*

(Received 30 September 1994)

Methods to estimate the number of degrees of freedom of chaotic dynamical systems suffer from intrinsic errors. The errors are due to the finite extent of phase space and are felt at any finite number of phase space points. We compare the errors of two methods to extract dimensions from scaling properties. One is based on the scaling of the number of points in spheres with varying radius and the other one concerns the scaling of the radius of spheres that contain a varying number of points. We argue that the latter method is preferable and we derive an analytic expression for the error. We compare both this systematic error and the error due to statistical fluctuations in different realizations of random sets to the results of numerical simulations.

PACS number(s): 05.45.+b, 47.53.+n

### I. INTRODUCTION

Temporal disorder in physical systems can be caused by the nonlinear interplay of just a few degrees of freedom. The estimate of a lower bound of their number from a signal produced by the system is a worthwhile goal. If this number indeed turns out to be small, one could try to determine the unstable periodic orbits of the system, try to influence the system by stabilizing these orbits, try to predict future states of the system, or even try to find a model in the form of nonlinear ordinary differential equations. The estimate of the number of degrees of freedom, therefore, is a first step towards a better understanding of the source of disorder.

The dimension estimates that we will discuss in this paper are based on scaling arguments. For example, the well known Grassberger-Procaccia correlation integral [1] is the scaling of the number  $C(r)$  of phase space points in balls of radius  $r$ ,  $C(r) \sim r^D$ , where  $D$  is the dimension of the phase space. An alternative way, the practical implementation of which predates the Grassberger-Procaccia algorithm (GPA), is based on the scaling of the radius  $r$  of balls that contain a given number of points  $k$ ,  $r(k) \sim k^{1/D}$  [2]. The first method is called a fixed size method, the second a fixed mass method.

For a given number  $N$  of phase space points, the fluctuations of the scaling function  $C(r)$  increase with decreasing  $r$ . This is because the filling of phase space becomes increasingly sparse at smaller distances. On the other hand, at large distances the boundary of phase space is felt. Scaling behavior, therefore, is restricted to an interval bounded by these extremes. The problem is that the scaling region shrinks rapidly with increasing dimension of phase space.

Crudely, the nearest neighbor distance in  $D$  dimensions is  $\delta_1 = N^{-1/D}$ , and the average distance to the boundary of a  $D$ -dimensional hypercube is  $\delta_2 = 1/(2D + 2)$ . With increasing  $D$  the size of the scaling interval for a given number of points  $N$  shrinks to zero if  $\delta_1 = \delta_2$ , or

$$N^{-1/D} = \frac{1}{2D + 2}, \quad N = (2D + 2)^D. \quad (1)$$

The number of phase space points where the scaling interval has just collapsed increases superexponentially with increasing dimension.

Smith [3] pointed out the restriction on the accessible range of dimensions. The simple Eq. (1) was refined by allowing for an error in the estimated dimension and requiring a finite scaling dynamic range. Nerenberg and Essex [4] argue that Smith's approach is too restrictive. Scaling may actually extend to distances smaller than  $\delta_1$  and distances larger than  $\delta_2$ . For the correlation integral  $C(r)$ , the small-distance limit is then given by the maximum allowable statistical fluctuations of  $C(r)$ .

As was realized in [4], the effect of the proximity of the phase space boundary is not a sharp cutoff at large scales but is felt at *all* values of  $r$ . Therefore, the effect of the boundary proximity is to introduce a deviation from scaling that, unlike the expression for  $\delta_2$ , depends on  $N$ .

Compared to the correlation integral, the scaling of fixed mass methods extends to the smallest mean nearest neighbor distance of the set and is not a compromise involving the size of statistical errors. For fixed mass methods, the lower bound is close to  $\delta_1$  but scaling extends to much larger distances. It is therefore a worthwhile goal to estimate the boundary error for this method. However, because the distance is now the dependent variable, the effect of the proximity to the boundary is much harder to estimate.

Geometric limits on scaling are but one problem in detecting the presence of low-dimensional chaos in experimental data. We have designed an expression for the boundary effect for data that consists of white noise that is uniformly distributed in hypercubes. However, it is well known that the correlation integral for colored noise can lead to spuriously small dimensions [5].

The scaling of asymptotic orbits of dynamical systems concerns their organization in phase space. A point in

phase space represents the projection of the system on (maybe linear) modes of motion. In theory, a faithful representation of the instantaneous state of a  $D$ -dimensional system with at least  $D$  independently measured projections could be bypassed by embedding a measured time series [6]. The proper choice of embedding parameters presents yet another problem in an estimate of  $D$  from a time series [7]. We will not dwell upon the question of embedding because we believe that well instrumented experiments in physics can actually measure spatially distributed information that can be put to use to more effectively reconstruct the underlying phase space.

In Sec. II we will review the analysis of Nerenberg and Essex [4]. In Sec. III we give a brief derivation of the near-neighbor (NN) method and we examine the geometrical effects that lead to systematic errors. For an analytical calculation we have been forced to make some approximations. We check our analysis in Sec. III by comparing its results with those of numerical simulations involving random white noise.

## II. CORRELATION INTEGRAL

The correlation dimension  $d_c$  derives from the scaling with  $r$  of the average number of points  $C(r)$  in spheres with radius  $r$ ,  $C(r) \sim r^{d_c}$  [1]. Practically, one fixes the radius  $r$  and measures the mass contained in those spheres. Therefore,  $r$  is the independent variable and  $C$  is the dependent one. This observation serves to distinguish the correlation integral (GPA) from fixed mass methods, where the role of independent and dependent variables is exchanged.

The correlation integral  $C(r)$  is defined as

$$C(r) = \lim_{N \rightarrow \infty} \frac{2}{N^2} \sum_{i < j=1}^N \Theta(r - |\mathbf{x}_i - \mathbf{x}_j|), \quad (2)$$

where  $\Theta(r)$  is the Heaviside step function and the  $\mathbf{x}_i$ 's are state vectors in the system's phase space. Figure 1 shows the correlation integral for  $N = 10^4$  points that are distributed randomly in a  $D = 6$  dimensional hypercubic space. It is seen that before the slope of  $C(r)$  in a log-log plot starts approaching its nominal value  $d_c = 6$  at small distances, the size of its fluctuations has increased catastrophically.

The range of distances  $r$  over which scaling can be observed is bounded from above by the proximity of the phase space boundary. The reason is that for radii  $r$  that are comparable to the total size of phase space  $R$ , the correlation integral  $C(r)$  no longer increases with  $r$ . For a finite number of phase space points  $N$ , the value of  $C(r)$  at *small* distances is strongly fluctuating due to the sparseness of points at small scales. Therefore, the behavior of  $C(r)$  at small distances cannot be used to extract a dimension. Based on the idea that a dimension derived from scaling acquires a systematic error when these two regions have started to overlap, Smith [3] derived a criterion for the minimum number of points  $N_{\min}$

needed to be able to conclude a dimension  $D$  within 5%,  $N_{\min} = 42^D$ .

In [4] an attempt is made to quantify more accurately the error in the correlation dimension of finite size data sets. The key point is the shape of the correlation function in the presence of boundaries. An equation was derived for  $C(r)$  in the case of  $D$ -dimensional hypercubic spaces that are uniformly filled with points. Using this formula, an estimate was made for the dimension underestimation for a given  $N$  and  $D$ . The dimension error has two contributions: one,  $\Delta_b d_c$ , due to the boundary proximity and one,  $\Delta_s d_c$ , due to the statistical fluctuations of the value of the correlation integral.

The systematic part of the dimension error  $\Delta_b d_c$  can be made smaller by moving the interval  $[r_0, \xi r_0]$  over which the slope of  $C(r)$  is determined in a log-log plot to smaller  $r_0$ . However,  $C(r)$  at the smallest values of  $r$  is most affected by statistical fluctuations and the contribution  $\Delta_s d_c$  to the error in  $d_c$  will be largest. A compromise between the two types of error then leads to an optimal

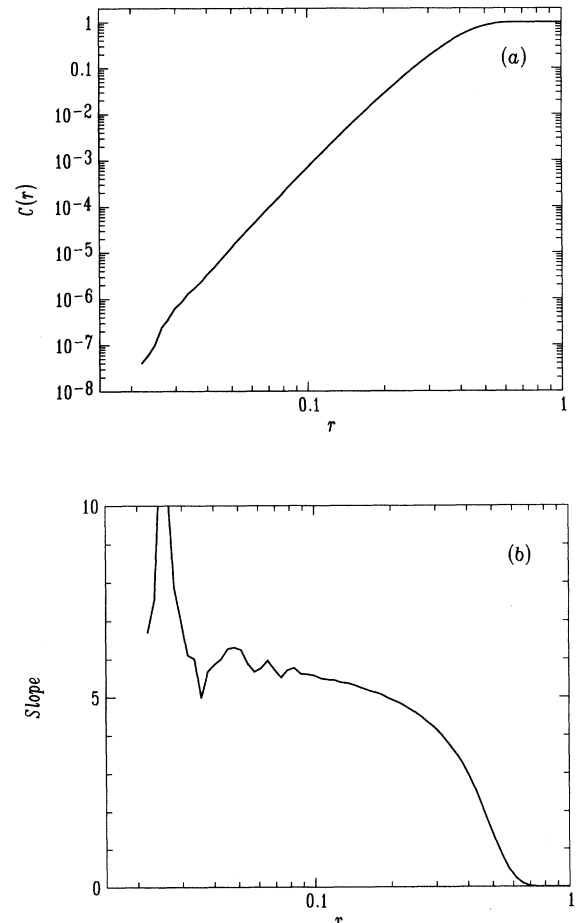


FIG. 1. (a) Correlation integral  $C(r)$  for  $N = 10^4$  points distributed randomly in a hypercubic space with dimension  $D = 6$ . (b) Local slope of scaling curve of (a).

choice of the interval  $[r_0, \xi r_0]$ .

In the next section we will demonstrate that no such compromise is needed for the fixed mass method, and we will derive an expression for the dimension error in this complementary method.

### III. NEAR-NEIGHBOR METHOD

The near-neighbor (NN) method for estimating the dimension of a set is concerned with the scaling of the mean near-neighbor distance  $\delta$  as a function of the size  $n$  of the set and as a function of the neighbor order  $k$  (nearest neighbor is  $k = 1$ , next-nearest neighbor is  $k = 2$ , and so on). Crudely,

$$\delta(n, k) \simeq n^{-1/D} k^{1/D}. \quad (3)$$

A more precise analysis involves important finite-size corrections that play an explicit role in the resulting scaling functions. These corrections are important for small values of  $k$ . Also, the fixed radius method has finite size corrections, but those disappear for the correlation integral (which is but one version of the fixed radius method) [8].

Let us assume that we have a set of  $N$  points that are distributed uniformly in a  $D$ -dimensional phase space. Consider a subset of  $n$  phase space points from the original set and a reference point  $i$ . The probability to find  $k$  elements of this subset within a radius  $r$  of the reference point  $i$  is

$$S_i(r; k, n) = \binom{n}{k} p_i^k (1 - p_i)^{n-k}. \quad (4)$$

Here  $p_i$  is the probability to find one phase space point of the subset in a sphere of radius  $r$  around  $i$ . This probability is proportional to the ratio of the volume of the sphere and the volume of phase space  $p_i = \rho K_D r^D$  with  $\rho = (V_D R^D)^{-1}$ , and  $K_D = (\pi)^{D/2} / \Gamma(1 + D/2)$  the volume of the unit sphere. The geometrical factor  $V_D$  is  $V_D = 2^D$  for a hypercubic space and  $V_D = K_D$  for a hyperspherical space. It is important to notice that near each of the reference points  $i$ , we have assumed the probability  $p_i$  to scale as  $p_i \simeq r^D$ . Therefore, we ignore the issue of multifractality which we believe is a moot point when trying to estimate the nearest integer value of the dimension of a large-dimensional attractor.

The probability to find the  $k$ th near-neighbor with a distance between  $r$  and  $r + dr$  from the point  $i$  is the probability to find  $k - 1$  elements within the sphere around the point  $i$  times the probability to find the  $k$ th element of the set in the spherical shell  $[r, r + dr]$ . This last probability is proportional to the volume of the hyperspherical shell times  $n$ .

$$\begin{aligned} P_i(r; k, n) dr &= S_i(r; k - 1, n) n \rho D K_D r^{D-1} dr \\ &= n \rho D K_D r^{D-1} \\ &\quad \times \frac{(n \rho K_D r^D)^{k-1} \exp(-n \rho K_D r^D)}{\Gamma(k)} dr, \end{aligned} \quad (5)$$

where we have used the Poisson approximation for  $S_i(r; k - 1, n)$ , for large values of  $n$ ,  $n \gg k$ , such that  $(1 - p_i)^{n-k} \approx \exp(-n p_i)$  and  $\binom{n}{k-1} \approx n^{k-1} / \Gamma(k)$ . Averaged over different realizations of the distribution of  $n$  points over the attractor, the value of the  $k$ th near-neighbor distance of reference point  $i$  is

$$\begin{aligned} r_i(k, n) &= \int_0^\infty dr r P_i(r; k, n) \\ &= \frac{\Gamma(k + \frac{1}{D})}{\Gamma(k)} (n \rho K_D)^{-1/D}. \end{aligned} \quad (6)$$

For a single reference point  $i$ , Eq. (6) defines the point-wise dimension. A more adequate definition will more faithfully sample the attractor, and an average over reference points needs to be done [9]. Here we will use a simple linear average. Averaging  $r_i(k, n)$  over  $M$  reference points defines  $r(k, n) = \frac{1}{M} \sum_{i=1}^M r_i(k, n)$ . For large

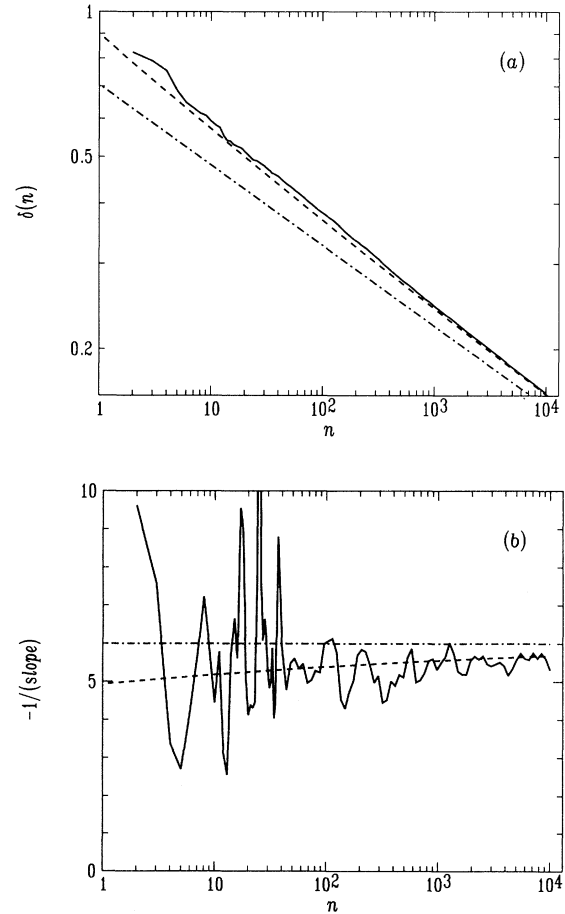


FIG. 2. (a) Full line: nearest neighbor distance  $\delta(n)$  for  $N = 10^4$  points distributed randomly in a hypercubic space with dimension  $D = 6$ . Dashed line: prediction of Eq. (11). Dash-dotted line: prediction of Eq. (11) without correction for the effect of the boundary proximity. (b) Local dimensions that follow from scaling curves in (a), where the slope is the local slope of the scaling curves in (a).

$k$  we have  $\Gamma(k + \frac{1}{D})/\Gamma(k) \sim k^{1/D}$  and we regain Eq. (3). Scaling behavior can be found by measuring the  $n$  dependence of  $r(k, n)$  at fixed  $k$  or by studying its  $k$  dependence at fixed  $n$ . As the latter involves finite size corrections for finite  $k$  [the factor  $\Gamma(k + 1/D)/\Gamma(k)$ ], we will concentrate on the first. Incidentally, these finite size corrections are trivially accounted for in a fitting procedure that is used to derive  $d_n$  from a log-log plot of  $r(k, n)$  vs  $k$  [10].

More general averages, that are averages of  $r_i$  raised to a certain power,  $r(k, n) = \left(\frac{1}{M} \sum_{i=1}^M r_i^\gamma(k, n)\right)^{1/\gamma}$ , lead to a whole spectrum of dimensions  $D^q$  [10,11]. For the correlation integral  $q$  takes the value  $q = 2$ . For the scaling of nearest neighbor distances that are averaged with  $\gamma = 1$ , the value of  $q$  depends on the dimension as  $q = 1 - 1/D$ . Effectively, therefore, the dimension that is estimated with the near-neighbor method is close to the information dimension when  $D$  is large. It can be shown that for all choices of  $\gamma$ , the value  $q = 2$  is an upper limit [10].

Practically, the NN method is implemented by first selecting randomly a set of reference points. Next, a random subset of  $n$  points is selected from the total number of  $N$  phase space points and the near-neighbor distances to each member of the set of reference points is computed. The number  $n$  is (exponentially) increased and the process is repeated. Of course, when the size  $n$  of the subset approaches the total number of points  $N$ , subsequent random sets are no longer independent. One would expect that these dependencies affect the scaling; however, we have found this effect to be negligible. We have designed an analytical model for it whose results are consistent with those of numerical simulations.

Figure 2 shows a scaling curve of the nearest neighbor distance  $\delta(n) \equiv r(k = 1, n)$  for a uniform random distribution of  $10^4$  points in a six-dimensional unit cube. Clearly, the slope of the scaling curve for the largest  $n$  tends to the nominal value  $-1/6$ , but is everywhere smaller. The corresponding value of the estimated dimension, therefore, is always smaller than 6. The scaling curve of Fig. 2 is complementary to the correlation integral: the slope most closely approaches  $-1/6$  where its statistical fluctuations are *smallest*. Because the average is taken over distances, the smallest distance of Fig. 2 is much larger than the smallest  $r$  of Fig. 1, which is the pair distance of the single, most dense spot of the attractor. On the other hand, the scaling of the near-neighbor distance extends to much larger  $r$ . Unlike for the correlation integral, there is *no ambiguity* as to the relevant scaling range of near-neighbor distance curves; it is  $[N/\xi, N]$ , where  $\xi$  is the scaling dynamical range.

### A. Dimension error

To understand the reason of the underestimation of the dimension, we notice that the volume of hyperspheres that intersect the boundary of phase space is reduced. Crudely, the nearest neighbor distance is given by the requirement that  $nC(r) \approx 1$ , where  $C(r)$  is the probability to find one point in a sphere with radius  $r$ . When  $n$  is small,  $r$  is large and the chances that such a sphere inter-

sects a boundary are large. To compensate for the volume reduction due to the intersection,  $r$  becomes larger compared to the situation where no boundaries are present. Conversely, when  $n$  is large, the probability for intersection is drastically reduced, and  $r$  approaches the unbounded situation. Therefore, the boundary effect causes the slope of the scaling curve  $\ln \delta(n)$  versus  $\ln n$  to be more negative, resulting in an underestimate of the dimension.

Naively, the fraction of reference points that see the boundary decreases with decreasing  $r$  as  $r^{D-1}$ . Therefore, it is in principle possible for fixed radius methods to exclude from the average at given  $r$  those reference points that are closer to the boundary than  $r$ . Because  $r$  in the near-neighbor method is the *dependent* variable, no such separation in points whose apparent neighborhood scaling is not affected by the boundary proximity is possible. This circumstance also makes an analytical estimate of the boundary effect for the NN method much harder. Such an analytical estimate is precisely what is attempted here.

From now on we will consider nearest neighbors ( $k = 1$ ) only and we will accordingly drop the  $k$  dependence. The effect of the boundary is that the function  $r_i(n)$  depends on the location of the reference point  $i$  with respect to the boundary. Because we will consider intersections with a single bounding hyperplane only,  $r_i(n)$  only depends on the distance  $l$  of the reference point  $i$  to its nearest phase space bounding hyperplane. Therefore, we define a function  $\tilde{r}(l; n) \equiv r_i(n)$ . The phase space average then takes the form

$$\delta(n) = \int_0^R dl g(l) \tilde{r}(l; n), \quad (7)$$

where the geometric structure factor  $g(l)dl$  expresses the probability to find a reference point that has a distance  $l'$  to the nearest boundary, with  $l'$  in the interval  $l' \in [l, l + dl]$ . For homogeneously and isotropically randomly filled hypercubic or hyperspheric spaces with linear size  $R$  this is

$$g(l) = \frac{D}{R} \left(1 - \frac{l}{R}\right)^{D-1}. \quad (8)$$

The function  $\tilde{r}(l; n)$  is the near-neighbor distance of reference points that are at a distance  $l$  from the boundary. The  $l$  dependence of this function merely expresses that the scaling for reference points near the boundary will be different from points in the interior of the phase space volume. For a reference point at a distance  $l$  from the boundary, the probability to find a nearest neighbor at distance  $r$  in a set of  $n$  points is a slight generalization of Eq. (5),

$$P(r; l; n) = n\rho \frac{\partial V(r; l)}{\partial r} \exp[-n\rho V(r; l)], \quad (9)$$

where  $V(r; l)$  is the volume of a sphere of radius  $r$  with its center at a distance  $l$  from the boundary. For  $r < l$  the sphere does not intersect the boundary and we have  $V(r; l) = K_D r^D$ ; for  $r > l$  the volume  $V(r; l)$  is the volume of a chopped hypersphere.

For the calculation we take for the phase space a  $D$ -dimensional randomly filled hypercube with edge length  $2R$  or a hypersphere with radius  $R$ . The hypercube is bounded by  $2D$  different hyperplanes; the hypersphere by a single spherical boundary.

For a reference point that is at a distance  $l$  from the boundary, the expression for the ensemble-averaged near-neighbor distance  $\tilde{r}(l; n)$  separates into two terms, one where the presence of the boundary is not yet felt, and one that involves the volume of a chopped sphere around the reference point:

$$\tilde{r}(l; n) = \int_0^l dr r P(r; l; n) + \int_l^{2R-l} dr r P(r; l; n). \quad (10)$$

Note that the upper integration limit  $2R - l$  of the second integral is given by the restriction that in the case of hypercubical spaces we consider intersections with a single boundary only. In general, multiple intersections lead to the necessity of evaluating  $(2D + 1)$ -fold integrals for  $\tilde{r}(l; n)$ , in which case calculations would no longer be tractable. The error resulting from our approximation of the volume will be largest for large  $r$  and large  $D$ , where we overestimate the volume of the hypersphere around reference point  $i$ . However, because of the exponential behavior of  $P(r; l; n)$ , the effect on  $\tilde{r}(l; n)$  for large  $r$  in the integration will be small. The effect of our approximation of the chopped volume will, however, become sizable for any  $r$  at large  $D$ . This is because in high-dimensional phase spaces most points are near the edge of phase space. Finally, for a hyperspheric space we will neglect the curvature of the boundary for  $r > l$ . This again will result in an overestimate of  $V(l; r)$ , especially for large  $r$  and large  $D$ . The given expression for  $\tilde{r}(l; n)$  needs to be averaged over phase space using the structure factor  $g(l)$ . Both the computation of  $\tilde{r}(l; n)$  and the averaging are described in the Appendix with the simple result Eq. (A11),

$$\delta(n) = n^{-1/D} A_1(D) + n^{-2/D} A_2(D), \quad (11)$$

with

$$A_1(D) = R \Gamma\left(1 + \frac{1}{D}\right) \left(\frac{V_D}{K_D}\right)^{1/D},$$

$$A_2(D) = R \Gamma\left(\frac{2}{D}\right) \left[ \left(\frac{V_D}{K_{D-1}}\right)^{2/D} \Lambda(D) - \left(\frac{V_D}{K_D}\right)^{2/D} \right], \quad (12)$$

and where the function  $\Lambda(D)$  is explicated in Eq. (A12). The effect of the boundary is given by the second term of Eq. (11). Because it is  $O(n^{-2/D})$  and the regular scaling is  $O(n^{-1/D})$ , the relative effect of the boundary on the scaling function vanishes as  $n^{-1/D}$  for large  $n$ .

From Eq. (11) we compute the dimension  $d_n$  by fitting a straight line in a plot of  $\ln \delta(n)$  versus  $\ln n$  over an assumed scaling interval  $[N/\xi, N]$ . Obviously, however, Eq. (11) does not display simple scaling behavior, but our choice is motivated by the standard practice of analyzing experimental data.

$$d_n = -\frac{\ln(N) - \ln(N/\xi)}{\ln \delta(N) - \ln \delta(N/\xi)}. \quad (13)$$

The boundary effect leads to a systematic error  $\Delta_b d_n = d_n - D$  in this dimension estimate,

$$\begin{aligned} \Delta_b d_n &\approx D^2 \frac{A_2(D)}{A_1(D)} N^{-1/D} \frac{1 - \xi^{1/D}}{\ln \xi} \\ &\approx -DN^{-1/D} \frac{A_2(D)}{A_1(D)}, \end{aligned} \quad (14)$$

where we have assumed that  $|N^{-1/D} A_2(D)/A_1(D)| \ll 1$ . In the Appendix we will show that  $A_2(D)$  is positive; therefore, the error has *negative* sign. In agreement with our qualitative arguments the proximity of the boundary leads to an *underestimate* of the dimension.

Apart from the systematic error on the dimension due to the boundary proximity, there exists a statistical error due to the fluctuations of the near-neighbor distances in different realizations of  $D$ -dimensional random sets. Unlike for the fixed radius method (the correlation integral), however, the choice of the scaling interval does *not* depend on this statistical error. It is precisely this circumstance that renders application of fixed mass (near-neighbor) methods less ambiguous than dimension estimates with fixed radius methods. The scaling interval is simply located at the smallest possible distances (largest possible  $n$ ), where the statistical error and the error due to the boundary proximity are smallest.

An estimate of the statistical error in  $\delta(n)$  due to sample-to-sample fluctuations can be computed from the distribution function of near-neighbor distances  $P_i(r; k, n)$ , Eq. (5),

$$\Delta \delta_i = \left[ \int_0^\infty r^2 P_i(r; k, n) dr - \left( \int_0^\infty r P_i(r; k, n) dr \right)^2 \right]^{1/2}. \quad (15)$$

We assume that these fluctuations are uncorrelated for different reference points. The statistical contribution  $\Delta_s d_n$  to the dimension error then is

$$\begin{aligned} \Delta_s d_n &= \frac{2D^2}{\ln \xi} N^{-1/2} \frac{1}{\Gamma(1 + 1/D)} \\ &\quad \times [\Gamma(1 + 2/D) - \Gamma^2(1 + 1/D)]^{1/2}, \end{aligned} \quad (16)$$

where the factor  $N^{-1/2}$  is the reduction of the statistical error due to the average over  $N$  reference points.

## B. Numerical results

We have used numerical simulations to ascertain our estimates of systematic dimension errors. In these simulations we have randomly filled  $D$ -dimensional hypercubes and hyperspheres with  $N$  points. All  $N$  phase space points are used as reference points. No average was done over different random subsets. Such an average will only work at small sizes  $n$  and serves to reduce

the fluctuations in the scaling function such as shown in Fig. 2(b). However, the dimension is estimated using the large- $n$  behavior of the scaling function. Averaging over subsets at large  $n$  is not effective because different subsets are then no longer independent.

In Fig. 2 the analytic expression for  $\delta(n)$  [Eq. (11)] is compared to the result of a numerical simulation. It is seen that the addition of the term  $O(n^{-2/D})$  not only results in a better approximation of the local slope, but also in a better approximation of the absolute size of  $\delta(n)$ . Of seven realizations of random sets with  $1 \leq D \leq 10$  and  $N = 10^4$  we have measured the apparent dimensions  $d_n$  and the error  $\Delta d_n = d_n - D$ . It is important to notice that we have not used embedding of a random time series in  $D$  dimensions using the method of delays [6].

In Fig. 3 the results of simulations on hypercubes are compared with the prediction of Eq. (14). Our analytic estimate of the boundary proximity effect is for values of  $D$  up to 6, in good agreement with the result of the simulation. The apparent dimension of our simulated phase space fluctuates from realization to realization. The size of the fluctuations decreases with increasing  $\xi$ , i.e., with enforcing a larger dynamical range. The result for an analogous simulation but now for hyperspherical spaces is shown in Fig. 4. Clearly, our analytic expression Eq. (14) performs better for hypercubical spaces. For other choices of  $N$  and  $\xi$  the results are the same.

In the case of phase spaces with a hyperspherical boundary there does not exist the problem of finding out with which of the bounding surfaces the neighborhood sphere intersects and it is possible to improve upon our analytic formula. The result of a numerical evaluation of Eq. (A3) is also shown in Fig. 4. It is in better agreement

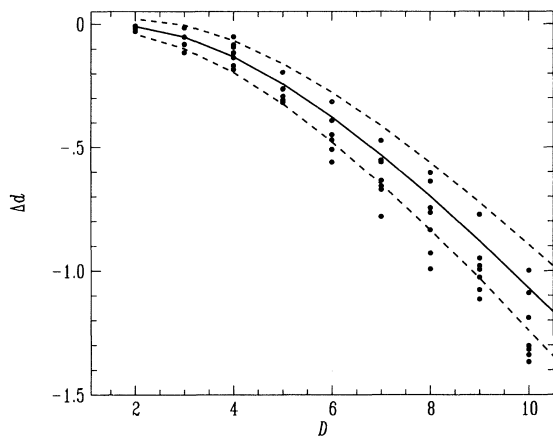


FIG. 3. Dimension error of the near-neighbor method. Dots: results of numerical simulations for the error made when estimating the dimension of randomly filled hypercubes with dimension  $D$ . The number of points is  $N = 10^4$ ; the scaling dynamical range is  $\xi = 2$ . Solid line: prediction of Eq. (14) with  $V_D = 2^D$ . Dashed lines:  $\Delta_b d_n + \Delta_s d_n$  and  $\Delta_b d_n - \Delta_s d_n$ , respectively, with the statistical error  $\Delta_s d_n$  given by Eq. (16).

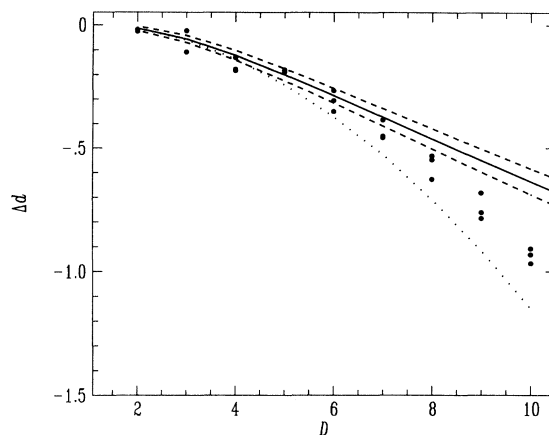


FIG. 4. Dimension error of the near-neighbor method. Dots: results of numerical simulations for the error made when estimating the dimension of randomly filled hyperspheres with dimension  $D$ . The number of points is  $N = 10^4$ ; the scaling dynamical range is  $\xi = 10$ . Solid line: prediction  $\Delta_b d_n$  of Eq. (14) with  $V_D = K_D$ . Dashed lines:  $\Delta_b d_n + \Delta_s d_n$  and  $\Delta_b d_n - \Delta_s d_n$ , respectively, with the statistical error  $\Delta_s d_n$  given by Eq. (16). Dotted line: boundary error computed from numerically evaluating Eq. (A3).

with the results of the simulation and demonstrates that Eq. (14) may be performing poorly for large dimensions  $D$ .

Figure 5 shows the results of simulations for the correlation integral. The underestimate of the dimension is slightly less than for the near-neighbor method, but, as explained, the GPA method suffers from an ambiguous choice of the scaling interval. The choice made in [4] of

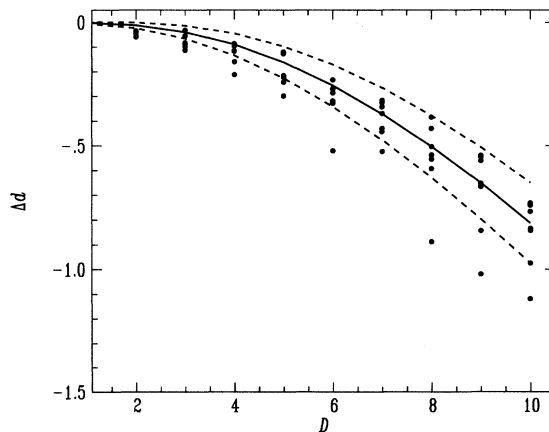


FIG. 5. Dimension error of the correlation integral. Dots: results of numerical simulations for the error made when estimating the dimension of randomly filled hypercubes with dimension  $D$ . The number of points is  $N = 10^4$ ; the scaling dynamical range is  $\xi = 2$ . Line: prediction of [4].

the lower bound of the scaling interval has made the size  $\Delta_s d_c$  of the fluctuations of the dimension estimate comparable to that of Fig. 3. Shifting it to smaller  $r$  results in an increase of  $\Delta_s d_c$ .

When comparing the results of Figs. 3 and 5 it should be realized that straight lines were fitted to the scaling curves. As the abscissa of the fixed mass scaling function in Fig. 2(a) spans a much larger dynamical range than that of the fixed radius method in Fig. 1(a), the *relative* dynamical ranges used in the dimension estimates are very different.

#### IV. CONCLUSION

The underestimation of dimensions is a serious flaw of scaling methods for dimension measurements. Due to this effect, one may be tempted to conclude a small number of degrees of freedom when this number is actually so large that it eludes measurement.

The underestimation is partly caused by geometric effects. We have shown that the fixed radius (the correlation integral) and fixed mass methods (the near-neighbor method) suffer from this problem to approximately the same degree. However, the correlation integral has the additional problem of the ambiguity in the choice of the scaling interval. No such ambiguity is present for the fixed mass method, where the scaling is always determined by the smallest distances where both the fluctuations of the scaling function and the boundary effect are smallest.

We summarize our result in Fig. 6 by displaying the systematic dimension error [Eq. (14)] for several values of  $N$ . Assume an unknown dynamical system that has produced  $N$  samples of a  $D$ -dimensional trajectory. If for this set of  $N$  points an apparent dimension  $d'$  is found with  $\Delta d < d' - D < 0$ , it falls within the systematic error

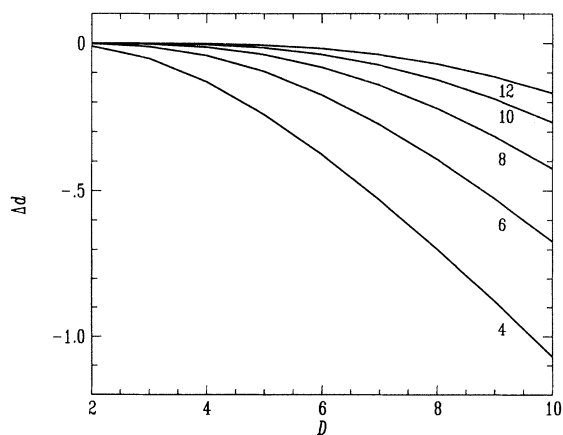


FIG. 6. Underestimation of dimensions using the near-neighbor method in case of  $D$  dimensional randomly filled hypercubes. The scaling range is  $\xi = 2$ ; the number of phase space points varies from  $N = 10^4$  to  $N = 10^{12}$ .

bound of the dimension estimate and it follows that the measured trajectory cannot be distinguished from one produced by space filling white noise. It is a striking observation that the convergence of  $\Delta_b d$  to zero with increasing  $N$  is so slow; hardly anything is gained when increasing  $N$  from  $10^{10}$  to  $10^{12}$ .

It is also a striking and counterintuitive observation that decisions about large dimensions can be made on basis of a few phase space points. For example, Fig. 6 suggests that it is possible to distinguish  $D = 8$  from  $D = 10$  using only  $10^4$  phase space points. From applying these methods of dimension estimate to experimental data we have learned that such a small error may be too optimistic. It illustrates the necessity of considering other sources of error in dimension estimates. A well documented source of error is the effect of time correlations of the phase space signal. It has been studied extensively in the context of the correlation integral [5].

The effect of the boundary on dimension estimates depends on the shape of the boundary. From our simulations it appears that the effect is slightly smaller in spherically bounded spaces than in spaces with a cubical boundary. Large-dimensional systems are in general systems that explore many degrees of freedom, i.e., systems that are described by partial differential equations. Strings of coupled nonlinear maps of length  $L$  are now widely accepted as faithful models of spatiotemporal chaotic dynamics. For such systems the dimension grows in proportion to their size; here  $D \sim L$ . It has been suggested [12] that the chaotic attractors of coupled map lattices for finite  $L$  may have zero thickness in some directions of  $L$ -dimensional phase space. It is clear that the existence of such internal boundaries would aggravate the problem of dimension estimates using scaling methods.

An interesting recent suggestion [13] has been to divide out the effect of boundaries on dimension estimates by normalizing the correlation integral on its value for uniform noise in the given phase space. For the near-neighbor method such a normalization may be done by fitting measured scaling curves with Eq. (11). If this method appears viable, dimension estimates may remain a valuable tool for analysis of spatially extended nonlinear systems. However, as the present paper again shows, they should be applied with great caution.

#### ACKNOWLEDGMENTS

It is a pleasure to acknowledge discussions with Niek Lopes Cardozo and Chris Schüller. This work was performed as part of the research programme of the association agreement between the “Stichting voor Fundamenteel Onderzoek der Materie” (FOM) and Euratom with financial support from the “Nederlandse Organisatie voor Wetenschappelijk Onderzoek” (NWO) and Euratom.

#### APPENDIX

In this appendix we will calculate the effect of the boundary on the scaling of the nearest neighbor distance

$\delta(n)$ , where  $n$  is the size of a random subset of the total set with  $N$  elements. The average nearest neighbor distance for a subset of size  $n$  is

$$\delta(n) = \int_0^R dl g(l) \int_0^l dr r P(r; l; n) + \int_0^R dl g(l) \int_l^{2R-l} dr r P(r; l; n), \quad (\text{A1})$$

with the structure factor

$$g(l) = \frac{D}{R} \left(1 - \frac{l}{R}\right)^{D-1}. \quad (\text{A2})$$

The effect of the proximity of the boundary is contained both in the structure factor  $g(l)$  and in the volume of hyperspheres  $V(r; l)$  that enters the definition of  $P(r; l; n)$ . In order to compute the boundary effect in hyperspherical spaces we first have to estimate the volume of a sphere that is at a distance  $l$  from the spherical boundary of a  $D$ -dimensional phase space. This computation is illustrated in Fig. 7. In the case that  $r < l$ , there is no intersection and  $V(r; l) = K_D r^D$ ; for the case  $r > l$  we have

$$V(r; l) = K_{D-1} r^D \int_{-1}^{(l-h)/r} dy (1-y^2)^{(D-1)/2} + K_{D-1} R^D \int_{(R-h)/R}^1 dy (1-y^2)^{(D-1)/2}, \quad (\text{A3})$$

with  $h = (r^2 - l^2)/[2(R-l)]$  (see Fig. 7). If  $r \ll R$  and  $l \ll R$ , we have  $h \approx 0$  and the intersection of a sphere with radius  $r$  with the spherically bounded phase space with radius  $R$  is a plane. In this approximation the curvature of phase space is neglected and finally gives

$$V(r; l) = \begin{cases} K_D r^D, & r < l \\ K_{D-1} r^D \int_{-1}^{l/r} dy (1-y^2)^{(D-1)/2}, & r \geq l \end{cases} \quad (\text{A4})$$

From Eq. (A3) and Fig. 7, we see that  $2R-l$  is the

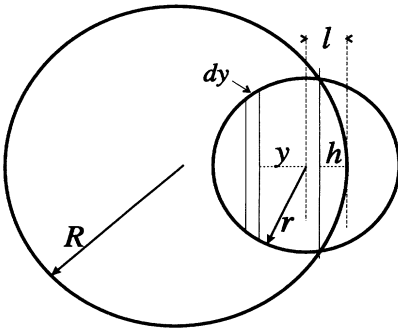


FIG. 7. Geometry for the calculation of  $V(r; l)$  for  $r > l$  in the case of a hyperspheric phase space.

maximum value of the radius  $r$ .

In the case of hypercubic spaces in principle we have to take into account intersections with more than one bounding plane. This would require the calculation of  $(2D+1)$ -fold integrals. We will use the approximation that spheres intersect with a single bounding hyperplane only. In our approximation, therefore, hypercubic and hyperspherical phase spaces are treated similarly. The only distinction is the substitution of the appropriate geometrical factor  $V_D$  that equals  $2^D$  in case of hypercubic phase spaces and  $K_D$  for hyperspherical spaces.

In the case of nearest neighbors the distribution function  $P(r; l; n)$  can be written as

$$P(r; l; n) = -\frac{\partial}{\partial r} \exp[-n\rho V(r; l)], \quad (\text{A5})$$

and performing partial integrations in Eq. (A1) gives

$$\delta(n) = -\int_0^R dl g(l) (2R-l) \exp[-n\rho V(2R-l; l)] + \int_0^R dl g(l) \int_0^l dr \exp[-nsV(r; l)] + \int_0^R dl g(l) \int_l^{2R-l} dr \exp[-nsV(r; l)]. \quad (\text{A6})$$

Because  $V(2R-l; l) = O(R^D)$ , the volume of phase space, the first term is of order  $e^{-n}$  and will accordingly be neglected. The computation of the remaining two terms of Eq. (A6) is simplified if we write the structure factor

$$g(l) = -\frac{dG(l)}{dl}, \quad \text{with } G(l) = \left(1 - \frac{l}{R}\right)^D. \quad (\text{A7})$$

The second term of Eq. (A6) involves the volume  $V(r; l)$  of spheres that are not intersected by a boundary. Using partial integration,

$$\int_0^R dl g(l) \int_0^l dr \exp[-n\rho V(r; l)] = \int_0^R dl G(l) \exp(-n\rho K_D l^D). \quad (\text{A8})$$

Approximating  $G(l) \approx (1 - Dl/R)$  and extending the upper integration limit to infinity, we at once recognize the emergence of two  $\Gamma$  functions, one multiplying the ordinary scaling  $n^{-1/D}$  and the other one associated with the boundary effect  $n^{-2/D}$ .

The third term of Eq. (A6) involves the volume of a sphere that is chopped by a bounding hyperplane. The



double integral can be simplified considerably by taking  $R$  instead of  $2R - l$  as the upper limit of the integration over  $r$ . This is justified because the neglected part becomes exponentially small with increasing  $n$ . The double integral over the region  $[0, R; l, R]$  can be done by introducing a new integration variable  $a = l/r$ .

$$\begin{aligned} & \int_0^R dl g(l) \int_l^{2R-l} dr \exp[-n\rho V(r; l)] \\ &= \int_0^R dr r \int_0^1 da g(ar) \exp[-n\rho K_{D-1} r^D f(a)], \end{aligned} \quad (\text{A9})$$

where

$$f(a) = \int_{-1}^a dy (1 - y^2)^{(D-1)/2}. \quad (\text{A10})$$

When the integration over  $r$  is extended to infinity, a power series in  $n^{-1/D}$  results with  $\Gamma$  functions as coefficients. The lowest order term is  $n^{-2/D}$  and involves the zeroth order term of the power series expansion of  $g(l)$ . Collecting all terms up to order  $n^{-2/D}$  we finally have

$$\begin{aligned} \delta(n) = R & \left\{ n^{-1/D} \Gamma\left(1 + \frac{1}{D}\right) \left(\frac{V_D}{K_D}\right)^{1/D} \right. \\ & + n^{-2/D} \Gamma\left(\frac{2}{D}\right) \left[ \left(\frac{V_D}{K_{D-1}}\right)^{2/D} \Lambda(D) \right. \\ & \left. \left. - \left(\frac{V_D}{K_D}\right)^{2/D} \right] \right\} + O\left(n^{-\frac{3}{D}}\right), \end{aligned} \quad (\text{A11})$$

with

$$\Lambda(D) = \int_0^1 da f(a)^{-2/D}, \quad (\text{A12})$$

which has been calculated numerically for different values of  $D$ . Using the fact that  $K_D/K_{D-1} = f(1)$ , it is trivial to show that the term between square brackets in Eq. (A11) is positive. As argued in Sec. III A, this implies that the boundary proximity causes dimensions to be underestimated. We note that the only way that the shape of the phase space volume (hyperspherical or hypercubical) enters in Eq. (A11) is through the geometrical factor  $V_D$ . This is because of our assumption that a sphere cuts a (much) smaller sphere as a plane. However, for other geometries one would have to consider different functions  $g(l)$ .

- 
- [1] P. Grassberger and I. Procaccia, Phys. Rev. Lett. **50**, 346 (1983); Physica D **9**, 189 (1983).
  - [2] K. W. Pettis, T. A. Bailey, A. K. Jain, and R. C. Dubes, IEEE Trans. Pattern Anal. Mach. Intell., **PAMI-1**, 25 (1979).
  - [3] L. A. Smith, Phys. Lett. A **133**, 283 (1988).
  - [4] M. A. H. Nerenberg and C. Essex, Phys. Rev. A **42**, 7065 (1990).
  - [5] A. Provenzale, L. A. Smith, R. Vio, and G. Murante, Physica D **58**, 31 (1992).
  - [6] F. Takens, in *Dynamical Systems and Turbulence*, Lecture Notes in Mathematics Vol. 898, edited by D. A. Rand and L. S. Young (Springer-Verlag, New York, 1981), p. 366.
  - [7] M. Ding, C. Grebogi, E. Ott, T. Sauer, and J. Yorke, Physica D **69**, 404 (1993).
  - [8] P. Grassberger, Phys. Lett. A **128**, 369 (1988).
  - [9] Strictly speaking, the order of averaging over realizations and averaging over reference points (local scalings) is wrong. The correct order is crucial when studying multifractal sets that have negative dimensions.
  - [10] W. van de Water and P. Schram, Phys. Rev. A **37**, 3118 (1988).
  - [11] R. Badii and A. Politi, Phys. Rev. Lett. **52**, 1661 (1984); J. Stat. Phys. **40**, 725 (1985).
  - [12] A. Torcini, A. Politi, G. P. Puccioni, and G. D'Alessandro, Physica D **53**, 85 (1991).
  - [13] M. Bauer, H. Heng, and W. Martienssen, Phys. Rev. Lett. **71**, 521 (1993).