

Even when data are defined in a high-dimensional space, they commonly lie (exactly or approximately) on a hypersurface of much lower intrinsic dimension (ID). Characterizing the ID of the data may be relevant within the preprocessing of the data, preliminary to more refined analyses: for instance, the data may be projected within a lower-dimensional space with some dimensionality-reduction technique, setting the dimension of the target space to the estimated ID. In addition, the ID may provide *per se* relevant information about the structure of the data.

Among techniques to estimate the ID, a large class of methods is based on the fact that, for sufficiently random data, nearest-neighbor (NN) distances follow scaling laws that depend on the ID. In particular, by assuming that the density of points is locally constant, one retrieves scaling laws that depend *solely* on the ID (i.e., no other parameters are involved), which can be leveraged for ID estimation. Various ID estimators were proposed, based on slightly varying choices of the scaling function and inference algorithm. Differences concern the number k of nearest neighbors considered, the ID-dependent scaling function used, the precise procedure to infer the ID from the scaling law, and the possible inclusion of additional information beyond distances [Ceruti et al., 2014]. ID estimation algorithms encounter three main difficulties undermining their accuracy and robustness: i) boundary effects, ii) variations of the density of points in the data, iii) undersampling in high dimension.

In their manuscript entitled “manifold-adaptive estimation revisited” the Authors propose an updated version of one the ID estimation algorithm proposed by Farahmand et al. 2007 (FSA estimator). Let us denote by r_k the ratio of the distances of k -th and the $2k$ -th NNs of a point, and let D be the ID value: if the density is locally constant, one approximately has $r_k = 2^{-D}$, or $D = -1/\log_2(r_k)$. Farahmand proposed to collect the $r_k(x)$ for each point x in the data, and estimated D as the mean of mode of $D(x) = -1/\log_2(r_k(x))$. The Authors build on the same idea, but find the exact distribution of $r_k(x)$ and consequently of the “local estimates” $D(x)$, under the assumption of locally constant density (i.e., the density can be assumed to be constant within a $2k$ -neighborhood of each point). Finding the exact distribution of $D(x)$, they show that the median of $D(x)$ (rather than the mean or mode) is equal to D . Importantly, this holds independently of k , a fact which allows selecting small values of k (even $k = 1$) for estimation. This median-based procedure, called mFSA estimator, yields an improvement over the estimation method proposed by Farahmand et al. 2007. A further improvement is given by a boundary-effect correction: the Authors estimate D on uniformly-sampled hypercubes, showing that boundary effects induce a systematic underestimation of D , with a relative error scaling exponentially with D . They then introduce a correction compensating for this systematic effect, which eventually gives the corrected median FSA estimator (cmFSA). They show that cmFSA is equally or more accurate than current state-of-art estimation methods, at least on test data proposed by Campadelli et al. 2015. Finally, they show that cmFSA can be applied to intracranial recordings to detect signal abnormalities occurring during epileptic seizures.

The manuscript is sound and well written. References are generally exhaustive. The methodology is clearly explained.

cmFSA seems to be competitive with state-of-art methods, and I believe it offers some advantages with respect to some of the above mentioned issues of ID estimators (in particular, i) and ii)). I think the manuscript is a fair contribution to the field of ID estimation, and it can be of interest to researchers in this area, and more in general to researchers needing accurate ID estimation as part of their data analysis pipelines. Therefore, I recommend the paper for publication in PeerJ.

However, I think that some improvements are needed before the paper is published. **Major points:**

- In my opinion, the main problem with the boundary-effect correction is that it is optimized for uniformly-sampled hypercubes, and may lead to overestimation of the ID in cases when the data are not uniformly sampled. This is clearly visible from table I: while the estimation is nearly perfect for uniformly sampled data on linear subspaces [M2,M9,M10a-c], or generally uniformly sampled data on locally flat spaces [M5,M7,M13], it yields an overestimation in the case of non-uniformities, such as the Gaussian case [M12], the non-linear manifold case [M6], or the sphere [M1]. The overestimation may be even more severe for non-uniform samplings with heavy-tailed distributions, such as the Cauchy distribution used in Facco et al. 2015. The authors should extensively comment on this point.
- Since this is a methodological work, I would recommend that the authors make publicly available the code implementing cmFSA.
- It is not clear how the different sample sizes were included in the calibration of the correction term. It seems that the calibration term used to infer the ID of the datasets M1-M13 was inferred from the $n = 2500$ hypercubes. Is one going to use the same term with datasets of different n ? It seems that one should rather use a term calibrated on that specific n . The authors should comment on this point. Furthermore, why was $k = 5$ used for calibration, instead of $k = 1$ used in subsequent analyses?

Minor points:

- The Authors may better stress the fact that their median-based procedure is independent of k , and thus allows selecting a minimal neighborhood size ($k = 1$). In this case, the used statistics is essentially equivalent to the one used by Facco et al., 2017 - even though the estimation procedure is slightly different. As in Facco et al., using a minimal size neighborhood can make the method very robust to density variations and curvature.
- The simplicity of the proposed statistic makes it suitable to be embedded within mixture-based approaches to provide better ID estimates when the ID is varying in the data set (Haro, G., Randall, G. & Sapiro, G. Translated poisson mixture model for stratification learning. Int. J. Comput.

Vis. 80, 358–374 (2008); Allegra M, Facco E, Denti F, Laio A, Mira A (2020) Data segmentation based on the local intrinsic dimension. Sci Rep 10(1):16449).

- The Authors may better clarify Eqs. (1-2). In Eqs (1-2), k is used to indicate both a variable quantity and a fixed quantity. In eq. (1), k is a variable quantity, like R [notice that Eq. (1) now uses both R and R_k , inconsistently]. In eq. (2), k is a fixed value, like r_k and r_{2k} . Also, the quantities in Eq. (1) should be better defined. I would recommend something like: “A usually basic assumption of kNN ID estimators is that the fraction of points f in a spherical neighborhood centered at x is approximately determined by the intrinsic dimensionality (D) and radius (R) times a – locally almost constant – mostly density-dependent factor ($\eta(x, R)$):

$$f/n = \eta(x, R)R^D$$

[...] If R_k is the distance at which the k -th neighbor is found, from Eq. (1) one can take the logarithm...”

- In Eq. (5), the Authors may better clarify what $p(r|k, K-1, D)$ is: something like “the probability that the normalized distance of the k -th neighbor among K neighbors is r if the intrinsic dimension is D .”
- “Thus, we can compute the pdf of the estimated values as plugging in $K = 2k$ into Eq. 5 followed by change of variables”(p. 4). This sentence might be more clearly rephrased, e.g., “Combining (5) and (6), one can obtain the pdf of the FSA estimator”
- In theorem 1, the Authors may mention that the substitution $a = 2^{-D/d_k}$ is monotonic, which justifies the invariance of the median.
- “This means that the median of the FSA estimator is equal to the intrinsic dimension independent of neighborhood size”. Again, this fact should be stressed because it allows using small k , which cannot be done in standard FSA: indeed, for small k , as evident from fig. 1, the mean and mode produce severe underestimates.
- In Fig. 2, the Authors may add a third panel showing on a simple plot the standard error of the median as a function of $\log(n)$, for different values of D (different curves for different values of D).
- I would put a derivation of Eq. (17) in the SI. (the rationale of the binomial is nearly obvious, but a full explanation may help the reader).
- Are periodic boundary conditions used in Figure 4, as the main text indicates? This should be clarified also in the caption of Fig. 4, to stress the difference with Fig. 3, which is not using PBC.

- In eqs. (21)-(22) it would be better to bring in some notational clarity. What are d, D, \hat{d} ? Note that D was always used as the true value of intrinsic dimension.
- How is the error in Fig. 6 defined? It is stated that “the error rate – the fraction of cases, when the estimator did not find (missed) the true dimensionality”. What does this mean exactly? That $|D_j - d_{ij}| > 1$?
- In fig. 7, what are 1-8 on the x axis? Is it simply the electrode number? To what areas do the grid recordings Gr-A ... Gr-F correspond? The Authors should specify it in Methods, or at least provide a reference.
- Why was $k = 10$ used in the analysis of electrode data? If results change a lot between $k = 1$ and $k = 10$, it may be because data were not optimally subsampled.
- In Methods, p. 13: Fig. 7 \rightarrow Sfig 2
- In fig. S2, I would stress that panel c shows that the error distribution after correction is approximately Gaussian. * “we observed a diagonal gradient of intrinsic dimensions on the cortical grid (Gr)” (p. 7). It is difficult to interpret a diagonal gradient (as opposed to a vertical gradient representing cortical hierarchy).
- there are some typing errors/language mistakes:
 - while resting state and during epileptic seizures (p. 1) \rightarrow between normal resting state and during epileptic seizures
 - MIND_ML (p. 2) \rightarrow MIND_KL,
 - Kullback-Leibner \rightarrow Kullback-Leibler (p. 2),
 - ((Levina and Bickel, 2015)) \rightarrow (Levina and Bickel, 2015)
 - cmFSA and DANCo was evaluated \rightarrow cmFSA and DANCo were evaluated (p. 6)
 - froto-basal \rightarrow fronto-basal (p. 7)
 - chose \rightarrow choose (SI p. 1)
 - distance, we assume \rightarrow distance, if we assume (SI p. 2)