

Associate Editor: I agreed with both reviewers that recommend Major Revision of the manuscript. However, these intensive comments implied that substantial revisions were needed.

Reviewer #1

A potentially useful contribution, albeit very descriptive. However, in its current form there are too many issues with the analyses for it to be publishable. The English expression also needs substantial work.

== General comments.

1. The hexagons used are not equal area and are therefore not comparable spatial units. The analyses are therefore incorrect and need to be redone.

Indeed, we realized this and have completely redone the analysis using a Cylindrical Equal Area projection. This now bases our analyses on equal area grids. A description of the procedure were included in section 2.2.1 *Spatial Representativeness Analysis*; page 7.

2. It is possible the polygon data have also been incorrectly reported (Table 1), although I am unable to locate the geographic polygon data plotted in Fig 1 to check this (it appears not to be available from Costello *et al.*, 2017).

Thanks for this comment, we have now revisited the number of polygons, names and identification of them. Table 1 shows this revision and coincides with Costello *et al.* (2017)'s bioregions (see <https://doi.org/10.1038/s41467-017-01121-2>). We have also corrected the reference in the text and added to the manuscript ref list.

3. The SRI index is subject to issues when the denominator is low. Consider (1 / 2) vs (100 / 200) - these have the same ratio but the number of missing taxa is very different. The authors need to place the work in the context of Marcer *et al.*, (2022, <https://onlinelibrary.wiley.com/doi/full/10.1111/ecog.06025>) which is a much more thorough analysis of data, albeit not focusing on marine fish species.

Respuesta:

We appreciate this comment as it shows that we did not fully convey the meaning of the SRI index. While we recognize the need to account for the absolute gap, or missing spp, the objective of SRI is not associated with identifying the absolute magnitude of the information gap. Indeed, this index only indicates the relative magnitude of this gap. However, following your suggestions, we have decided to generate a new category for those cases where only one species per cell is recorded. In these situations the SRI index indicates that the cell is "adequately" sampled. To correct this error, we created the category "Insufficient Records" (see section 2.2.1 in page 7).

4. The Chao2 index requires incidences. How were these handled? I assume each record was treated as a unit of 1? Or were the hexagons used to define the incidences?

As explained in page 7, incidence were based on the record number in hexagonal cells.

5. It is unclear why the average of the richness estimators is used. It would perhaps be better to simply calculate the Chao2 (or Chao1) and then also take into account the

confidence intervals. And perhaps consider the ICE index instead as it uses more than just singletons and doubletons. See for example <https://github.com/AnneChao/SpadeR>

The recent literature reviewed here, on how to evaluate large data repositories for estimating species richness (see Mora et al., 2008; Troia & McManamay, 2017), show that different parametric and non-parametric methods usually generate different results in calculating species diversity, as these methods consider different attributes of the data (Walther & Morand, 1998; Gotelli & Colwell, 2001). One way to solve this is to use an estimator-averaging approach, which is argued to be statistically more accurate than a single model (Mora et al., 2008).

See additionally, the work by Gotelli & Colwell, 2001: (<https://doi.org/10.1046/j.1461-0248.2001.00230.x>) and Walther & Morand, 1998 (<https://doi.org/10.1017/S0031182097002230>) who consider an average of different diversity indexes as well.

6. The wrong reference has been cited for the marine bioregions. It should be Costello et al. (2017) Nat Comm, not Costello & Chaudhury (2017).

Respuesta: Thank you for this remark. We have indeed realized this error. We now cite the correct reference.

== Specific comments

7. P2, L18. What is a hands on concept?

We apologize for the lack of clarity here. We checked the manuscript by native speaker and added a clarification. We mean, concepts related to "species inventory", "taxonomic inventory" or "completeness of the inventory", terms used in contemporary studies to emphasize the need to record the diversity of species in a given ecosystem or for a given taxon. We have incorporated this clarification. (Page 3, 2nd paragraph)

8. P3, L14. Amon?

Respuesta: Thank you for this remark, we fixed this typo.

9. S2.2.1. These hexagons are not equal area if they are constrained to be 1x1 degrees. Indeed the maps in the supp mat show they are not equal area so therefore the analyses using the hexagons are not valid. As the meridians of longitude converge towards to the poles the true area of each hexagon approaches zero, so there are many more hexagons at latitude -80 compare with latitude -10. This needs to be corrected.

The projection of the grid is corrected to "cylindrical equal area", so that now the grid does have hexagonal cells of equal area. This methodological modification has been incorporated in section 2.2.1 Spatial Representativeness Analysis.

10. P4, L34. Why give the formula for the mean? Or are the authors trying to make the analysis look more complex than it really is?

We apologize that displaying the formula produces the wrong impression here. We just believe that a succinct description of the index may help some users to better understand the index.

11. P5. Just say the SRI is the ratio of observed to expected. And why not also the difference? Otherwise one suffers from the low denominator problem, e.g. is 1/2 as important as 100/200? (see also general comments).

The objective is only to know the degree of representativeness of the cell, that is, whether or not the observed richness corresponds to the expected richness, instead of how many species are missing to reach the expected richness levels. However, we have now considered this observation by identifying cells that have only one observation. We additionally note that SRI only considers cells that are sampled "adequately", as we have now created an additional category to evidence these cases.

12. P5, L6-12. Assignment of a continuous scale to such arbitrarily defined classes does not really help. Why is it that $SRI=0.599999$ is incomplete but 0.6000001 is complete, for example? Just use the observed values.

While somehow subjective, we do believe that, in this case, categories are more useful to broadly understand how the data is distributed. Using a continuous scale of values or simply the observed values could be overwhelming as there are large differences for biodiversity hotspot areas. Particularly in areas where richness indices are very low. The categorization intends to standardize the information in order to emphasize the bias in the data, particularly for certain geographic areas. This additionally highlights the large information gaps on a global scale.

13. P5, L16. "asses"

Respuesta: Thank you for this remark, we fixed this error.

14. S2.2.3. This is a long-winded description of a spatial intersection.

Respuesta: Thanks for your comment.

15. P6, L3. "Just about" -> "Approximately"

Respuesta: Thank you for this remark, we fixed this error.

16. Fig 1 caption. Describe the normalisation used. It is not given in Table 1. The bioregions are not defined in Costello & Chaudhary (see general comments).

The reference to the work of Costello et al., (2017) has been corrected. Panel "A" delivers the ecoregion number information that can be collated in Table 1. Table 1 has been corrected to deliver all information plotted in Figure 1: Species richness (panel B); family richness (panel C); Proportion of occurrence records by bioregion area (panel D); and Shannon index (panel E).

17. Table 1. Given the level of detail in Fig 1, the bioregion data are not sufficiently accurate to support a precision of 1 km². Report to the nearest thousand km². Please ensure the areas have been calculated using geodesic methods if the data have not been projected into an equal area coordinate system first.

Table 1, specifically the column "Area (km²)" has been corrected as suggested by the reviewer. In addition, the projection of the hexagonal grid to cylindrical equal area has also been corrected, and the values for each column have been updated.

18. Fig 2. The colour scheme used suggests a false dichotomy between the upper and lower classes. A continuous colour scheme should be used, e.g. from light to dark.

Thank you for this remark. The colors have been corrected following the suggestions of the second reviewer for a color palette that is visible to people with color blindness (Rochinni et al., 2023, <https://doi.org/10.1016/j.ecoinf.2023.102045>).

19. P6, L7. I don't know what is meant by this: "We used hexagonal grids that fit the geographic reality of marine ecosystems". Hexagonal grids are very useful because each hexagon's centroid is equidistant to its six nearest neighbours, unlike with rectangular tessellations. Hexagons also work better than square cells when tessellating data on a spheroid.

We mean that hexagons fit better by their geometric shape to the shape of marine bioregions, which are not exact squares. In other words, hexagonal cells fit much better to the shapefile boundaries of marine bioregions, which are not perfect polygons, they do not have right angles. The idea of using hexagonal cells is drawn from the work of Costello et al., (2017) (<https://doi.org/10.1038/s41467-017-01121-2>), here the possibility of delimiting bioregions using this grid shape is explored.

20. P12, L20. I am not sure of the relevance of this para. The example studies are presented as a shopping list without being integrated. Are all needed?

The purpose of this paragraph is to inform the reader about the various studies that have been carried out, which include different taxa and methodologies. The objective here is to contextualize and show that the systematic evaluation of large data repositories is a highly relevant topic, especially if the final objective is to generate inputs for biodiversity conservation. We have slightly corrected this paragraph to reflect this idea (Page 18).

21. P13, L14. This is incomplete. "the greater the overestimation of the index than richness". In any case, as the resolution becomes coarser one expects to see SR increase, or at least never decrease due to how it is calculated.

The larger the grid, the larger the area covered by the cell and the more species are, hence, considered. This, in turn, increases the value of SRI locally. So, using a lower grid resolution, will indicate that bioregions have better sampled areas as records may be accumulated in a very small area compared to the large size of a cell. Conversely, using a higher grid size resolution (e.g. 1°x1°) captures these very localized variations of the data. We have modified this paragraph so that this idea is properly conveyed (Page 19, line 19).

22. P13, L18. What is a thick grid? And why would hexagons also not suffer from this issue? Any regular tessellation will have this issue.

We refer to a low resolution grid. We have changed the word "thick" to "lower resolution". As we have already mentioned, a low resolution grid could be underestimating the SRI index by covering a large area when the records are clustered in more localized areas. This is especially true in coastal areas given the nature of the records hosted by GBIF and OBIS (opportunistic and citizen science records). On the other hand, we believe that a hexagonal shaped cell is better suited to the irregular shape of coastal areas, a methodology we adopt from Costello et al., (2017).

23. Fig 5. Why fit a straight line to what is clearly a curvilinear relationship? And if the dotted line is just to show a negative relationship then it is redundant as this is obvious. If significance is needed then see the first point in this comment.

Thank you, we have eliminated such line in the figures

24. P25, L7. This github repo is not available.

We apologize, it is now available.

Reviewer #2:

This is a well written and structured paper. I have some major points to be considered by the authors.

25. THE NEED FOR THIS STUDY: The pressing need for this study is not explained in detail or, at least, lacks a clear aim statement. Further, I would enlarge a bit the intro part on sampling effort which is, in my view, the most important spatial bias in species and diversity distribution mapping.

Respuesta:

While we understand the comments brought up by this reviewer, we believe to have clearly stated the context of the main concern in the introduction. The main idea behind this manuscript is that the large data accumulation of biodiversity information has not necessarily included an enlargement of our knowledge of biodiversity. In other words, we seem to have a fair amount of biased data providing a partial representation of our biota.

A large part of such discussion is related to developing a common framework to describe biodiversity in itself, and another part of the discussion is related to how we have sampled our biota from a macroecological perspective. This is extensively presented through the work of Yang et al (2013), Hortal et al (2008), Mora et al (2008), García-Rosello et al (2015) and Troia and McManamay (2017).

26. CARTOGRAMS: I strongly suggest in this paper the use of cartograms to represent in just one image both species diversity (as colour) and sampling effort (as a distortion of objects).

e.g.:

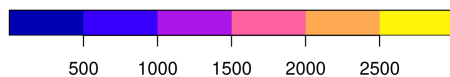
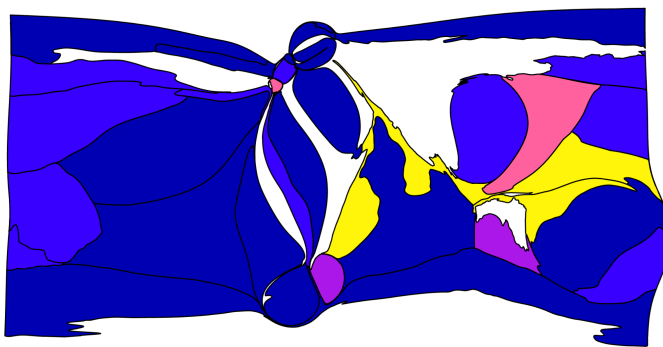
- Diffusion-based method for density-equalizing maps. Michael T. Gastner and M. E. J. Newman in PNAS

- Sampling biases shape our view of the natural world. Alice C. Hughes, Michael C. Orr, Keping Ma, Mark J. Costello, John Waller, Pieter Provoost, Qinmin Yang, Chaodong Zhu and Huijie Qiao in ECOGRAPHY

This would be one of the manners to represent spatial bias together with information on diversity. Showing uncertainty is mandatory in this type of paper.

See also the cartogram R package @: <https://CRAN.R-project.org/package=cartogram>

We tried to generate meaningful cartograms of Costello et al (2017)'s bioregion using species richness as colors and sampling effort as distortion. Unfortunately, we were not able to generate a representation of richness that will spark more information. We include here the best candidate for you to evaluate. The result is horrendous, and believe that this may be due to the fact that deformed marine bioregions completely shadows the shape of continents that help us understand the picture. We hence believe that we better stick to traditional representations.



27. Counts of records, species, families and Shannon diversity for each bioregion: they should be put as an infographic. The information provided as a table is too much. Translate Table 1 to an infographic or a graph and be put the complete table as supplementary material.

We have decided to keep the table as it allows us to evaluate the most succinct representation of the information used in this manuscript. It also provides lookup table to identify bioregions through their ID.

28. Spatial Representativeness Index: I definitively disagree with the colour ramp used to show the spatial representativeness.

See e.g.: <https://doi.org/10.1016/j.ecoinf.2023.102045> and change the colour ramp appropriately. Further, the table is complex and useless in my view.

All color palettes have been changed following the suggestions in Rocchini et al. (2023).

29. Protected area overlap and SRI: the same things I stated for the representation of counts apply here. This is also true for the overlapping FAO fishery exploitation areas and SRI grid.

We have decided to present our results on FAO fishery exploitation and protected areas in tables, as marine protected areas only represent an extremely small fraction of bioregions. Generating maps of this does not make much sense as it is poorly represented at a global scale. Regarding the FAO areas, the delimitations of these areas are not the same as the bioregions (they are totally different shapefiles), we believe that plotting this information on a map could be confusing for the reader. Our work already has two scales of analysis (bioregions and grids), we do not want to further complicate the information we are trying to communicate.

30. Records accumulation rate for each bioregion across the four decades analyzed: change red and blue with colours visible to everyone. This is also true for the map of slope values and for Spatial representativeness index (SRI) mapping of cells of different size.

All color palettes have been changed following the suggestions in Rocchini et al. (2023).

31. Why in the log relationships linear patterns are shown? I am clearly seeing a strong decay related to other types of functions to be used here.

Respuesta:

Thank you for this remark. We have now removed trend lines.

32. Compliments and thanks for the use of LaTeX which, as usual, renders everything more simple to read. I appreciated a lot the use of symbols for \ref to subsections, which is really straightforward.

Thank you very much for your comments and encouragement.

==Very minor:

33. Remove the final point in affil 4 after Chile.

Respuesta: Thank you for this remark, we fixed this typo.

34. Keywords in alphabetical order

Respuesta: Thank you for this remark, we fixed this error.