

Biases with the Generalized Euclidean Distance measure in disparity analyses with high levels of missing data

Lehmann, Oscar; Ezcurra, Martin; Butler, Richard; Lloyd, Graeme

DOI:

[10.1111/pala.12430](https://doi.org/10.1111/pala.12430)

License:

Other (please specify with Rights Statement)

Document Version

Peer reviewed version

Citation for published version (Harvard):

Lehmann, O, Ezcurra, M, Butler, R & Lloyd, G 2019, 'Biases with the Generalized Euclidean Distance measure in disparity analyses with high levels of missing data', *Palaeontology*, vol. 62, no. 5, pp. 837-849. <https://doi.org/10.1111/pala.12430>

[Link to publication on Research at Birmingham portal](#)

Publisher Rights Statement:

Checked for eligibility 04/02/2019

This is the peer reviewed version of the following article: Lehmann, O. E., Ezcurra, M. D., Butler, R. J. and Lloyd, G. T. (2019), Biases with the Generalized Euclidean Distance measure in disparity analyses with high levels of missing data. *Palaeontology*. doi:10.1111/pala.12430, which has been published in final form at: <https://doi.org/10.1111/pala.12430>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Use of Self-Archived Versions.

General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

Take down policy

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact UBIRA@lists.bham.ac.uk providing details and we will remove access to the work immediately and investigate.

1
2 BIASES WITH THE GENERALIZED EUCLIDEAN DISTANCE IN DISPARITY ANALYSES
3
4 WITH HIGH LEVELS OF MISSING DATA
5
6
7

8
9 by OSCAR E. R. LEHMANN^{1*}, MARTÍN D. EZCURRA^{1,2*}, RICHARD J. BUTLER², and
10
11 GRAEME T. LLOYD³
12
13
14

15
16 ¹ Sección Paleontología de Vertebrados, CONICET–Museo Argentino de Ciencias Naturales
17
18 “Bernardino Rivadavia”, C1405DJR, Buenos Aires, Argentina; e-mails: lehmannxii@gmail.com,
19
20 martindezcurra@yahoo.com.ar
21

22
23 ² School of Geography, Earth and Environmental Sciences, University of Birmingham, Edgbaston,
24
25 Birmingham B15 2TT, UK.
26

27
28 ³ School of Earth and Environment, University of Leeds, Leeds LS2 9JY, UK.
29

30
31 * Corresponding authors.
32
33

34
35 *Abstract.* The Generalized Euclidean Distance (GED) has been extensively used to conduct
36
37 morphological disparity analyses based on palaeontological matrices of discrete characters. This is in
38
39 part because some implementations allow the use of morphological matrices with high percentages
40
41 of missing data without needing to prune taxa for a subsequent ordination of the data set. Previous
42
43 studies have suggested that this way of using the GED may generate a bias in the resulting
44
45 morphospace, but a detailed study of this possible effect was still lacking. Here, we test if the
46
47 percentage of missing data for a taxon artificially influences its position in the morphospace, and if
48
49 missing data affects pre- and post-ordination disparity measures. We find that this use of the GED
50
51 creates a systematic bias, whereby taxa with higher percentages of missing data are placed closer to
52
53 the centre of the morphospace than those with more complete scorings. This bias extends into pre-
54
55 and post-ordination calculations of disparity measures and can lead to erroneous interpretations of
56
57 disparity patterns, especially if specimens present in a particular time interval or clade have distinct
58
59
60

1 proportions of missing information. We suggest that this implementation of the GED should be used
2 with caution, especially in cases with high percentages of missing data. Results recovered using an
3
4 alternative distance measure, Maximum Observed Rescaled Distance (MORD), are more robust to
5
6 missing data. As a consequence, we suggest that MORD is a more appropriate distance measure than
7
8
9
10
11 GED when analysing data sets with high amounts of missing data.
12
13
14

15
16 *Keywords.* Morphological disparity, distance measure, missing data, palaeontological matrices.
17
18
19

20 A large number of palaeontological studies dealing with morphological disparity have been published
21
22 over the last decade, mainly driven by the widespread use of already available morphological matrices
23
24 constructed for phylogenetic analyses. A critical step in most of these studies is the transformation of
25
26 the morphological matrix into a distance matrix. This is done using a distance measure, such as
27
28 Gower's coefficient (Gower 1971), the maximum observed rescaled distance (MORD; Lloyd 2016),
29
30 or the Generalized Euclidean Distance (GED; Wills 1998), among others. The existence of missing
31
32 data in the morphological matrix can produce distance matrices with missing entries when certain
33
34 distance measures are used. This means that no distance could be calculated for certain pairs of taxa
35
36 because of the absence of overlapping information between them, a rather common situation when
37
38 dealing with palaeontological information. A complete distance matrix is necessary if an ordination
39
40 is desired to reduce the dimensionality of the data set to produce a morphospace, calculate some
41
42 common disparity measures, or both. A strength of the GED is that it operates by inferring missing
43
44 distances and hence it is possible to always return a complete distance matrix, and in part because of
45
46 this it has become one of the most extensively used distance measures (*e.g.* Wills 1998; Brusatte *et*
47
48 *al.* 2011; Prentice *et al.* 2011; Thorne *et al.* 2011; Butler *et al.* 2012; Ruta *et al.* 2013a, b; Hetherington
49
50 *et al.* 2015; Oyston *et al.* 2015; Marx and Fordyce 2016; Lamsdell and Sendel 2017; Ezcurra and
51
52 Butler 2018). However, some concerns have been raised about the smoothing effect the GED may
53
54
55
56
57
58
59
60

1
2 have on the distance matrix and the disparity estimates derived from it, especially when it is used to
3
4 analyse matrices with a high percentage of missing data (Lloyd 2016, Ezcurra and Butler 2018).
5

6 However, some concerns have been raised about the smoothing effect the GED may have on
7
8 the distance matrix and the disparity estimates derived from it, especially when it is used to analyse
9
10 matrices with a high percentage of missing data (Lloyd 2016, Ezcurra and Butler 2018).
11
12

13 In this study we explore the performance of the GED in disparity analyses based on discrete
14
15 character matrices to address the concerns about its possible biases. This is especially important due
16
17 to the widespread use of the GED. The behaviour of an alternative distance measure, MORD, is also
18
19 analysed as a control against which to compare the GED. Three main objectives drove this study.
20
21 First, to determine if a hypothesised bias resulting from the replacements made by the GED are indeed
22
23 present; second, to explore if this bias produces a noticeable distortion in the morphospace; and finally,
24
25 to determine if differences in the percentage of missing data between groups of taxa can explain which
26
27 group is recovered as the most disparate.
28
29
30
31
32
33

34 **MATERIALS**

35
36
37
38

39 A data set of 158 published morphological matrices, each containing at least 20 taxa and 50 characters,
40
41 was used for this study. These matrices comprise a wide spectrum of taxonomic groups, including
42
43 plants, beetles, echinoids, actinopterygians, basal tetrapods, dinosaurs, birds, lepidosaurs, and
44
45 mammals, among others, and were primarily conceived for cladistic phylogenetic analyses. They
46
47 range from 20 to 223 taxa and from 50 to 486 characters. Only discrete characters are present. Most
48
49 of the matrices are compiled at the personal web page of one of the authors (GTL;
50
51 <http://www.graemetlloyd.com/>), were previously used by Wright et al. (2016), and all are available
52
53 in Lehmann *et al.* (2019, SI 2), together with a complete list of references. The analyses were
54
55 performed with different subsets of these matrices, as discussed below. Detailed information about
56
57 their dimensions, distribution of missing data, and use is available in Lehmann *et al.* (2019, SI 3-4).
58
59
60

METHODS

The following analyses were carried out using both the GED and the MORD for the generation of the distance matrices. In some cases, the ordination may become impossible when the MORD is used, as it is one of the distance measures that does not guarantee the creation of a complete distance matrix. To solve these cases, the TrimMorphDistMatrix function of Claddis (Lloyd 2016) was used to generate a complete distance matrix. This function iteratively removes the taxa responsible for the generation of most empty cells until the matrix is complete. These taxa were also trimmed from the GED matrix in order to generate comparisons with the same taxon sampling. This approach also makes the analyses presented here more conservative than if the GED were used directly on the untrimmed matrix, as it is commonly used, because the trimmed taxa tend to be the most incomplete, and thus the most prone to be affected by any bias generated by the GED.

All the analyses performed in this study were conducted in the programming environment R (R Core Team 2018), using custom-made functions and functions implemented in the package Claddis v. 0.2 (Lloyd 2016). The scripts used for this study are available in Lehmann *et al.* (2019). Statistical significance was assessed at $\alpha = 0.05$.

Calculation of GED and MORD

The GED (d_{ij}) between a taxon i and a taxon j is calculated with the following formula:

$$d_{ij} = \sqrt{\sum_{k=1}^v (S_{ijk}W_{ijk})^2}$$

where k is the current character, v is the total number of characters, S_{ijk} is the dissimilarity between taxon i and taxon j for character k , and W_{ijk} is the weight of character k . It can be used to generate a

complete distance matrix by replacing missing dissimilarities with a weighted mean value \hat{S} obtained from those dissimilarities that could be calculated, as follows:

$$\hat{S} = \frac{\sum_{k=1}^v (S_{ijk} W_{ijk} \delta_{ijk})}{\sum_{k=1}^v range(X_k) W_{ijk}}$$

where δ_{ijk} is 1 if both taxa have a scored, non-missing value for character k , and 0 if any of them is missing, and $range(X_k)$ represents the maximum possible dissimilarity for character k . Note that in Wills' (1998) original paper this weighted mean fractional univariate distance was intended to represent the taxon pair being considered only, but (as pointed out by Hopkins and St John 2018) Lloyd (2016) misinterpreted this as representing the mean for the whole matrix, i.e., all taxon pairs. Thus as Wills (1998) originally conceived the GED it would suffer from the same incomplete distance matrix problem as other metrics. However, here we use the common Lloyd (2016) implementation available in Claddis versions 0.1-0.2 as this is both the only currently offered implementation that will always return complete pairwise distances and the most frequently applied GED implementation in the recent literature.

The MORD is just a slight modification of Gower's Coefficient, yielding identical results to it when all the characters are binary or unordered (Lloyd 2016), and is calculated as:

$$d_{ij} = \frac{\sum_{k=1}^v (S_{ijk} W_{ijk})}{\sum_{k=1}^v range(X_k) W_{ijk}}$$

The distances obtained with this formula were subjected to the arcsine square-root transformation, as is default in the R package Claddis (Lloyd 2016). A detailed account of the two distance measures used here is available in Lehmann *et al.* (2019, SI 1).

Effect on the morphospace

The first effect evaluated was if the use of the GED influenced the position of taxa in the morphospace. This was done for an empirical set of morphological matrices and for a simulated set of matrices created by including sequentially higher percentages of missing data into empirical matrices with originally very low percentages (<5%) of missing data.

Empirical matrices with missing data. This analysis had the objective of testing if there is an inverse correlation between the percentage of missing data for a taxon and its distance to the centroid of the morphospace when the GED is used to construct the distance matrix. If this is the case, then taxa with a higher percentage of missing data will be closer to the centroid than those with fewer missing entries, arguably because the GED replaces a higher percentage of their pairwise dissimilarities with a weighted mean value. We further considered whether correlations may be more pronounced when smaller amounts of variance are used in the calculation of distance to the centroid and that the bias might disappear after a certain threshold level of explained variance.

A total of 126 matrices with more than 5% of missing data were used for this analysis. The first step was the creation of the distance matrices, which was followed by their ordination with a principal coordinates analysis (PCoA; Gower 1966; Legendre and Legendre 1998: ch. 9) to generate a morphospace of reduced dimensionality (Wills 2001). The Lingoes correction (Lingoes 1971) was applied if negative eigenvalues were created. This correction is necessary if all eigenvalues are required for the interpretation of the data (or if the absolute value of the largest eigenvalue is larger than the smallest positive eigenvalue of the principal coordinate axes [PCos] of interest [Legendre and Legendre 1998: p. 437]), but it frequently causes a distinct reduction in the variance explained by the first few PCos (Hopkins 2016). As the PCoAs are frequently corrected in published disparity analyses (*e.g.* Prentice *et al.* 2011; Butler *et al.* 2012; Oyston *et al.* 2015), we decided to include the correction as a step in the disparity pipeline, but we note that other authors prefer to avoid the correction if it is not completely necessary (*e.g.* Hopkins and Smith 2015; Hopkins 2016). Thus, to

1
2 test if the use of a correction could affect the results of our study, this particular analysis was also
3
4 performed without any correction for negative eigenvalues. PCos with negative eigenvalues were
5
6 excluded from the analysis.
7

8
9 Subsequently, the Euclidean distance was calculated between each taxon and the centroid of
10
11 the morphospace, which is its origin of coordinates and also the mean of the PCos. This was done
12
13 employing the PCos that correspond to accumulated variance from 5% to 100% in 5% increments,
14
15 and the particular case of the variance accumulated in the first three PCos, which are the most
16
17 commonly used to graphically represent the morphospace. Finally, the Spearman's correlation
18
19 coefficient was calculated between the percentage of missing data for each taxon and the
20
21 corresponding distance to the centroid of the morphospace. This was calculated for each distance
22
23 measure and each level of explained variance.
24
25
26
27
28

29 *Empirical matrices with simulated missing data.* To determine if the results obtained in the empirical
30
31 analyses could be simulated, some of the scored data entries in 33 matrices with less than 5% of
32
33 missing data were randomly replaced with missing values. These matrices have no possible
34
35 meaningful correlation between the percentage of missing data of their taxa and the distances of taxa
36
37 to the centroid. This is for either of two reasons: (1) most of the taxa had no missing data, which
38
39 caused the majority of the distances to the centroid to be evaluated based on a percentage of missing
40
41 data equal to zero; or (2) there were no missing entries in the whole matrix and, therefore, no standard
42
43 deviation (and thus no correlation coefficient) could be calculated.
44
45
46
47

48 To ensure a rather realistic distribution of missing values across the taxa, a beta distribution
49
50 with parameters $\alpha = 6$ and $\beta = 10.67$ was used to draw the proportion of missing data that each taxon
51
52 would have in the morphological matrix. The beta distribution has a domain of $[0, 1]$, and the selected
53
54 parameters ensured an expected value of 0.35, which is approximately the median proportion of
55
56 missing data observed in the empirical matrices used in the previous analysis. To avoid unrealistic
57
58 scenarios, such as a taxon or character with one or none scorings, the minimum number of non-
59
60

1 missing characters scored per taxon was limited to three and the minimum number of non-missing
2 scorings for a character to two. The disparity analysis protocol was the same as with the empirical
3 data set, with the distance to the centroid calculated using all the PCos. The entire protocol was
4 repeated 200 times on each of the 33 matrices.
5
6
7
8
9

10 11 12 13 *Effect on disparity measures*

14 We explored the effect that the GED may have on four commonly-used disparity measures: the pre-
15 ordination measures mean pairwise distance (MPD) and weighted mean pairwise distance (WMPD),
16 and the post-ordination measures sum of variances and sum of ranges. Details about the calculation
17 of these metrics and some of their interpretations can be found in Lehmann *et al.* (2019, SI 1). First,
18 we evaluated how the disparity values of a matrix varied as increasing percentages of entries were
19 randomly replaced with missing values. Secondly, we studied how the values of the disparity
20 measures for groups of taxa from the same matrix varied as the groups were assigned different
21 percentages of missing entries.
22
23
24
25
26
27
28
29
30
31
32
33

34
35
36 *Whole matrix analysis.* This analysis was designed to test if the use of the GED induces a bias
37 associated with the percentage of missing data of the morphological matrix into the chosen disparity
38 indices.
39
40
41
42

43 A total of 33 matrices with less than 5% of missing data were used for this analysis. The first
44 step was the calculation of the four disparity measures for the unmodified matrices in order to have a
45 baseline value to which compare the subsequent results. The next step involved the random
46 replacement of between 10% and 60% of the entries of the matrix with missing values, with
47 increments of 10%, and the calculation of the disparity measures for each case. The replacement with
48 missing entries started from the original, unmodified morphological matrix every time (a matrix with
49 60% of missing data was assigned all its missing values at once, not six times increasing by 10%).
50 This protocol was replicated 200 times for each level of percentage of missing data.
51
52
53
54
55
56
57
58
59
60

1
2 The MPD and the WMPD were calculated without trimming the distance matrix, but the
3
4 calculation of the sum of variances and the sum of ranges sometimes required the trimming of taxa,
5
6 as the PCoA was a necessary step for their calculation. As a result, the disparity estimates for these
7
8 post-ordination measures are comparable with each other but not with those of the pre-ordination
9
10 measures because they were calculated on potentially different taxon samplings. The sum of variances
11
12 and the sum of ranges were calculated using all the PCos.
13
14
15
16
17

18 *Groups of taxa from the same matrix.* This analysis was designed to evaluate if the disparity measures
19
20 calculated for groups of taxa within a matrix would be higher for groups with smaller percentages of
21
22 missing data, and lower for groups with higher percentages of missing data. If this is the case, the
23
24 effect may be attributed to the fact that the weighted mean value calculated by the GED would
25
26 generate a stronger bias in the groups with higher percentages of missing data, thus making them
27
28 more homogeneous when compared with groups with smaller influence from the replacement value.
29
30
31

32 To conduct this analysis, matrices that could be broken into five equally-sized, reasonably
33
34 large (10 or more taxa) groups were chosen. These matrices were also required to have less than 5%
35
36 of missing data. Seven of the previously used matrices with less than 5% of missing data fit these
37
38 criteria, and two additional matrices with an originally higher percentage of missing data but with
39
40 incomplete taxa removed until the threshold was met were also considered, summing up a total of
41
42 nine matrices. Five equally-sized groups of taxa were randomly chosen for each matrix. If the original
43
44 number of taxa in a matrix was not divisible by five, the maximum number of taxa that allowed
45
46 generating equally-sized groups was randomly chosen.
47
48
49

50 Thirteen different simulations of missing data distribution (Table 1) were conceived to provide
51
52 a diverse scenario of distributions of missing data among the groups of taxa. The simplest simulations
53
54 involved the random replacement of a given percentage of cells with missing data, equivalent to the
55
56 procedure explained in the previous section. The other simulations assigned different proportions of
57
58 missing data to each group of taxa. For example, Simulation 6 assigned 30% of missing data to Group
59
60

1
2 1, randomly distributed among all its taxa. Progressively smaller amounts of missing data, with steps
3
4 of 5%, were assigned to the other groups. This was designed to introduce a total of 20% of missing
5
6 data for the whole matrix, making Simulations 4–6 comparable because they all have the same overall
7
8 quantity of missing data and only differ in their distribution.
9

10
11 To assign the numeration to the groups, a procedure that ensured the independence between
12
13 the numeration and the percentage of missing entries introduced was implemented. The four disparity
14
15 measures were calculated for each group, and then the groups were ordered and numbered in a way
16
17 that made the Spearman's rank correlation coefficient between the disparity values for each measure
18
19 and the sequence 1 to 5 equal to 0. Then, 200 replicates were made for every combination of matrix
20
21 and simulation, and the disparity measures were calculated for each group in every case. The post-
22
23 ordination measures were calculated employing all the PCos.
24
25

26
27 The Spearman's rank correlation coefficient was calculated to compare the values of each
28
29 disparity measure and the group number. Because of the design of the simulations, the baseline
30
31 measures (*i.e.* with no simulation of missing data) and the group number have a correlation coefficient
32
33 of 0. In the simulations in which the distribution of missing data is heterogeneous between the groups,
34
35 the percentage of assigned missing data was the highest for Group 1 and sequentially lower, with the
36
37 lowest percentage in Group 5. If the disparity measures calculated for groups of taxa within a matrix
38
39 is indeed higher for groups with smaller percentages of missing data, and lower for groups with higher
40
41 percentages of missing data, the Spearman's correlation should be closer to 1, and equal to 1 in the
42
43 cases in which the ranking of the disparity measure is completely driven by the distribution of missing
44
45 data between the groups.
46
47
48
49
50
51

52 **RESULTS**

53 54 55 56 57 *Effect on the morphospace* 58 59 60

1
2 *Empirical matrices with missing data.* A total of 87.3% of the studied matrices showed a significant
3
4 and negative correlation between the distance to the centroid and the percentage of missing data in
5
6 the taxa when all the PCos were considered and the GED was used as the distance measure (Figures
7
8 1 and 2; see Lehmann *et al.* 2019, SI 5 for the complete results). This contrasts with the results
9
10 produced by the MORD, in which 79.4% of the matrices presented a non-significant correlation
11
12 between the variables. These percentages are similar to those obtained using only the first three PCos
13
14 (Figure 2). When using increasing amounts of variance (up to using all the PCos), the same result is
15
16 retained, *i.e.* mostly negative and significant correlations with the GED and mostly non-significant
17
18 correlations with the MORD (Figure 3; see Lehmann *et al.* 2019, SI 6 for the complete results). The
19
20 results of this analysis do not appear to be affected by the use of the Lingoes correction for negative
21
22 eigenvalues, as the percentages are similar when it is used and when no correction is applied
23
24 (Lehmann *et al.* 2019, SI 5).
25
26
27
28
29
30
31

32 *Empirical matrices with simulated missing data.* The simulations showed a trend similar to that
33
34 observed in the analysis with the matrices with more than 5% of missing data, in which the GED
35
36 produces a high proportion of significant correlations with negative coefficients, and the MORD
37
38 mostly non-significant correlations (Figure 4; see Lehmann *et al.* 2019, SI 7 for the complete results).
39
40
41
42

43 *Effect on disparity measures*

44
45 *Whole matrix analysis.* Both distances produced consistent results for the pre-ordination measures of
46
47 disparity (Figure 5). The GED generated the highest value for the complete matrix, with progressively
48
49 lower average disparities as the percentage of missing data increased. The MORD showed very
50
51 similar average disparity values across the range of percentage of missing data studied, with some
52
53 matrices having a small dip at around 60% of missing data (Lehmann *et al.* 2019, SI 8). This was
54
55 expected because the MORD explicitly takes into account the amount of missing data present.
56
57
58
59
60

1
2 The post-ordination measures presented a more complex pattern when the GED was used. The
3
4 sum of variances generally decreased as the percentage of missing data increased, starting from the
5
6 highest value if no additional missing data was introduced. However, this was not always the case. In
7
8 some of the matrices, the values of the sum of variances for the unmodified matrix were lower than
9
10 those with added missing data, while in others the maximum average value was reached with the
11
12 addition of 20–30% of missing data (Lehmann *et al.* 2019, SI 8). Nevertheless, the average sum of
13
14 variances tended to diminish after adding 30–40% of missing data. The sum of ranges presented
15
16 patterns similar to those of the sum of variances for each matrix.
17
18
19

20
21 When the MORD was used, the average sum of variances increased with higher percentages
22
23 of missing data in all the studied matrices. The sum of ranges presented a similar behaviour, but in
24
25 some cases the average disparity decreased during the addition of 10–30% of missing data and after
26
27 the introduction of higher percentages of missing data it increased.
28
29
30

31
32 *Groups of taxa from the same matrix.* When the GED was used, the nine studied matrices showed
33
34 correlations progressively closer to 1 as the difference in missing data between the groups increased
35
36 (Figure 6; Lehmann *et al.* 2019, SI 9 for the complete results). This effect was more pronounced for
37
38 the pre-ordination disparity measures, which also presented similar rankings of the disparity of the
39
40 groups when the missing data was homogeneously-distributed with respect to the ranking without
41
42 addition of missing data. The post-ordination measures showed a similar behaviour, but not as
43
44 pronounced. The average correlations with homogeneously-distributed missing data were negative,
45
46 though non-significant.
47
48
49

50
51 The simulations conducted with the MORD showed that the differences in the percentage of
52
53 missing data among the groups have little effect on the pre-ordination measures (Figure 6). For these
54
55 disparity measures, almost all of the mean correlation coefficients were not significantly different in
56
57 each level of percentage of missing data and no clear bias was present. In five of the nine studied
58
59 matrices the post-ordination measures progressively showed correlations closer to -1, which suggest
60

1
2 a reversed order of disparity from the base ranking. This implies that the groups with higher
3
4 percentages of missing data are recovered as more disparate than those with lower percentages, which
5
6 is the inverse pattern recovered for the GED.
7
8
9

10 **DISCUSSION**

11 *The effect of the GED in disparity analyses with missing data*

12
13
14
15
16 The results presented here reveal that the GED generally produces a significant negative correlation
17
18 between the percentage of missing data of a taxon and its distance to the centroid. In other words, the
19
20 more information is lacking for a taxon, the closer it will be placed to the centre of the morphospace
21
22 (Figure 7). This has been observed in the morphological matrices with a wide range of distributions
23
24 of missing data (Lehmann *et al.* 2019, SI 4), and in the almost complete morphological matrices with
25
26 the simulated inclusion of missing entries. Also, the presence of this bias seems to be independent of
27
28 the number of PCos used to calculate the distances to the centroid, so the use of only certain PCo axes
29
30 cannot solve this issue. In the analysis with groups with different percentages of missing data, the
31
32 GED showed a bias for all the studied matrices and for all the disparity measures explored. A result
33
34 consistent with this study has been recently reported by Ezcurra and Butler (2018) for an empirical
35
36 morphological matrix.
37
38
39
40
41
42

43 The bias caused by the use of the GED may have profound consequences for the study of
44
45 morphological disparity in the fossil record, as palaeontological data matrices, particularly those
46
47 focused on fossil vertebrates, typically have a high proportion of missing data (Lloyd 2016). If the
48
49 missing data are randomly distributed in the matrix, the effect will be diluted among all the taxa, but
50
51 if they are concentrated in some taxa (which is likely to be the case in many empirical data matrices),
52
53 the interpretations of the disparity analysis may be flawed. For example, if a clade is represented by
54
55 more incomplete fossil specimens in a particular time interval, the disparity measures calculated from
56
57
58
59
60

1
2 the GED could bias the results towards an interpretation in which this clade has a lower disparity in
3
4 this time interval than in others with more complete specimens.
5

6 In five of the nine studied matrices for groups with different percentages of missing data, the
7
8 MORD showed a bias towards recovering the groups with more missing data as more disparate in the
9
10 case of the post-ordination measures. This pattern is congruent with the observation of Ezcurra and
11
12 Butler (2018) that in an empirical data set some taxa with high percentages of missing data tended to
13
14 be pushed to the edge of the morphospace and away from the centroid when the MORD was used.
15
16 However, this bias can be mitigated by avoiding ordination altogether, as advocated by Lloyd (2016),
17
18 and as shown here in the pre-ordination (*i.e.* ordination-free) disparity measures.
19
20
21
22
23
24

25 *Contradictory results to previous studies*

26
27 Ciampaglio *et al.* (2001) studied the behaviour of the MPD and the sum of variances with different
28
29 amounts of missing data as part of the most thorough evaluation of the performance of several
30
31 disparity measures published so far. In that analysis, both measures showed no variation in their
32
33 average values as the percentage of randomly-distributed missing entries changed from 0 to 25%, by
34
35 increments of 5%. However, the results recovered in our study show a remarkably different behaviour,
36
37 with the mean values of the MPD decreasing as the percentage of missing data increases in the
38
39 simulations and the sum of variances showing a similar but more complex behaviour (Figure 5).
40
41
42

43 An important difference is how Ciampaglio *et al.* (2001) and we simulated the missing data.
44
45 Ciampaglio *et al.* (2001) simulated the missing entries by selecting iteratively a column of the matrix
46
47 (*i.e.* the scores of a character for all taxa) and then an entry of that column (*i.e.* a character score for
48
49 a single taxon). This selected entry was replaced by the mean of all the other scores in the column,
50
51 and the entire procedure was repeated until the desired proportion of “missing data” was achieved
52
53 (Ciampaglio *et al.* 2001: p. 698). However, this method is inadequate because in empirical matrices
54
55 the missing entries are not added sequentially. The successive replacement of missing entries with
56
57 the mean values of the column could never be applied to an empirical matrix, so the results yield no
58
59
60

1
2 information for comparisons with possible protocols of missing data handling. Also, it must be noted
3
4 that the matrices are effectively complete, as no scoring is actually missing.
5

6 In our study, the missing data is simulated by the direct elimination of information from the
7
8 morphological matrix. No attempt is made to replace the missing entries in the morphological matrix;
9
10 this is possible but outside the scope of this study. This leaves the handling of the missing information
11
12 to the calculation of the distance measure. The GED replaces the missing dissimilarities, not the
13
14 missing entries of the original data matrix; and the MORD does not replace any missing information,
15
16 it just rescales the distances with the number of available comparisons between a pair of taxa (see
17
18 Lehmann *et al.* 2019, SI 1). The GED then squares and sums the dissimilarities, and finally calculates
19
20 the square root of this value. This series of transformations complicates a straightforward prediction
21
22 of the behaviour of the calculated disparity measures when randomly-distributed missing entries are
23
24 added.
25
26
27
28

29 We replicated the procedure of missing data introduction of Ciampaglio *et al.* (2001) to
30
31 quantitatively compare their results to ours (see Lehmann *et al.* 2019, SI 10). Our results with the
32
33 protocol of Ciampaglio *et al.* (2001) are markedly different from those obtained by the original
34
35 authors. The relatively constant values of these disparity measures with increasing quantities of
36
37 missing data could not be recovered here, and we found that the absolute values of the disparity
38
39 measures tended to decrease as more missing data is introduced, resembling the results of our
40
41 simulations. For the sum of variances, the results were robust to either using or ignoring the Lingoes
42
43 correction. The calculation protocol of these measures also differs between both studies (see Lehmann
44
45 *et al.* 2019, SI 1), but the analysis we conducted showed that no matter which formula is used, the
46
47 results are consistent (see Lehmann *et al.* 2019, SI 10).
48
49
50
51

52 Thus, neither the different protocol of missing data introduction, the differences in the
53
54 calculation of the disparity measures, nor the use of a correction for negative eigenvalues (in the case
55
56 of the sum of variance) seem to account for the differences found.
57
58
59
60

1
2
3
4 *The MORD as a more robust distance measure for dealing with missing data*
5

6 Many methods have been explored to deal with large amounts of missing data (see Smith *et al.* 2014
7 and discussion therein). It is certainly difficult to imagine an elegant and unbiased way to do this,
8 especially because any imputation will most likely induce some kind of bias. In this study, we
9 investigated what happens when the GED is used to deal with missing data and found that the solution
10 it proposes induces an important bias in the pre- and post-ordination disparity measures and the
11 morphospace.
12
13
14
15
16
17
18
19

20 The GED was not designed to deal with matrices with large proportions of missing data, as
21 the use of a weighted mean dissimilarity was suggested for cases ‘Where the proportion of missing
22 data is relatively small [...]’ (Wills 1998: p. 471). The overall percentage of information being added
23 is higher than that of missing data found in the morphological matrix, because the GED replaces
24 missing character dissimilarities between pairs of taxa, not missing entries. For example, all the
25 matrices that had more than 31.5% of missing data in the original data matrix had more than 50% of
26 missing dissimilarities in the distance matrix. This fact implies that more artificial values were used
27 for the construction of their distance matrices than the information actually provided by the
28 morphological matrix. This does not mean that the distance measure itself yields inadequate results.
29 The main problem found in this work is that the weighted mean values added to permit the creation
30 of a complete distance matrix skews the results as was previously discussed, particularly when the
31 GED is applied to matrices with a large amount of missing data.
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

48 The use of alternative distance measures, such as the MORD, may alleviate the effect of
49 missing data. The MORD did not generate such a remarkable trend of correlations nor induce any
50 visible systematic bias in our analyses, at least in the majority of our results. If this is combined with
51 the suggestion that it retains more of the original signal of the data than other distances (Lloyd 2016),
52 the MORD seems to be a better alternative than the GED in the construction of the distance matrix
53 for a disparity analysis, particularly if the matrix being analysed has a high proportion of missing
54
55
56
57
58
59
60

1 data. An issue with the use of the MORD is that it does not guarantee a complete distance matrix,
2
3 thus precluding the ordination of the complete data set and the creation of a morphospace if missing
4
5 distances are present. In these cases, some taxa may have to be trimmed to continue with the analysis
6
7 or, as Lloyd (2016) recommended, ordination-free approaches can be pursued instead. Another
8
9 unexplored solution may be to combine the Wills (1998) and Lloyd (2016) interpretations of the GED
10
11 into a new “hybrid” GED such that missing distances are inferred from either the rest of the
12
13 information for the taxon pair (if available) or the global mean (all taxon pairs) if not. Thus a complete
14
15 pairwise distance matrix could still be returned, but the potential biases shown here could be
16
17 minimised.
18
19
20
21

22 The GED showed a systematic bias for the vast majority of the matrices studied, and the
23
24 MORD showed some hints of bias when coupled with the post-ordination disparity measures. From
25
26 these results, it is clear that, whichever distance is chosen, an adequate screening of the distance
27
28 matrices and the ordination results is extremely important to detect any possible bias, and avoid
29
30 artefactual conclusions, especially when post-ordination disparity measures are desired. The use of
31
32 the MORD may not solve the issue of incomplete distance matrices in these scenarios, but the
33
34 apparent solution provided by the GED seems to be unreliable not only for post-ordination disparity
35
36 measures, but also for some common pre-ordination ones. We thus suggest that the GED should be
37
38 used with caution, in particular when the morphological matrix has a high percentage of missing data.
39
40
41
42
43
44

45 *Acknowledgements.* We thank two anonymous reviewers for their useful comments about the
46
47 manuscript, and Sally Thomas for editorial suggestions. OERL thanks VLDS for proofreading the
48
49 first draft of this paper, and EDG for his help keeping the simulations running. RJB is supported by
50
51 the European Union's Horizon 2020 research and innovation programme under grant agreement
52
53 637483 (ERC Starting Grant TERRA).
54
55

56 **DATA ARCHIVING STATEMENT**

57
58

59 Data for this study, including data sets, scripts, and complete graphical results are openly available in
60 the Dryad Digital Repository: <https://datadryad.org/review?doi=doi:10.5061/dryad.4cv1421>

REFERENCES

- Baron, M. G., Norman, D. B. and Barrett, P. M. 2017. A new hypothesis of dinosaur relationships and early dinosaur evolution. *Nature*, **543**, 501-506.
- Brusatte S. L., Montanari, S., Yi, H. and Norell, M. A. 2011. Phylogenetic corrections for morphological disparity analysis: new methodology and case studies. *Paleobiology*, **37**, 1–22.
- Butler, R. J., Brusatte, S. L., Andres, B. and Benson, R. B. J. 2012. How do geological sampling biases affect studies of morphological evolution in deep time? A case study of pterosaur (Reptilia: Archosauria) disparity. *Evolution*, **66**, 147–162.
- Ciampaglio, C. N., Kemp, M. and McShea, D. W. 2001. Detecting changes in morphospace occupation patterns in the fossil record: characterization and analysis of measures of disparity. *Paleobiology*, **27**, 695–715.
- Ezcurra, M. D. and Butler, R. J. 2018. The rise of the ruling reptiles and ecosystem recovery from the Permian-Triassic mass extinction. *Proceeding of the Royal Society of London Series B, Biological Sciences*, published online 13 June 2018. doi:10.1098/rspb.2018.03.61.
- Gower, J. C. 1966. Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika*, **53**, 325–338.
- Gower, J. C. 1971. A general coefficient of similarity and some of its properties. *Biometrics*, **27**, 857–874.
- Hetherington, A. J., Sherratt, E., Ruta, M., Wilkinson, M., Deline, B. and Donoghue, P. C. J. 2015. Do cladistic and morphometric data capture common patterns of morphological disparity? *Palaeontology*, **58**, 393–399.
- Hopkins, M. J. and Smith, A. B. 2015. Dynamic evolutionary change in post-Paleozoic echinoids and the importance of scale when interpreting changes in rates of evolution. *Proceedings of the National Academy of Sciences of the United States of America*, **112**, 3758–3763, published online 23 February 2015. doi:10.1073/pnas.1418153112
- Hopkins, M. J. 2016. Magnitude versus direction of change and the contribution of macroevolutionary trends to morphological disparity. *Biological Journal of the Linnean Society*, **118**, 116–130.
- Hopkins, M. J. and St John, K. 2018. A new family of dissimilarity metrics for discrete character matrices that include inapplicable characters and its importance for, disparity studies. *Proceedings of the Royal Society B*, **285**, 20181784.
- Lamsdell, J. C. 2015. Horseshoe crab phylogeny and independent colonizations of fresh water: ecological invasion as a driver for morphological innovation. *Palaeontology*, **59**, 181-194.
- Lamsdell, J. C. and Selden, P. A. 2017. From success to persistence: identifying an evolutionary regime shift in the diverse Paleozoic aquatic arthropod group Eurypterida, driven by the Devonian biotic crisis. *Evolution*, **71**, 95–110.
- Legendre, P. and Legendre, L. 1998. *Numerical Ecology*. 2nd edition, Elsevier, Amsterdam.
- Lehmann, O. E. R., Ezcurra, M. D., Butler, R. J. and Lloyd, G. T. 2019. Biases with the generalized Euclidean distance in disparity analyses with high levels of missing data. Dryad Digital Repository. <https://datadryad.org/review?doi=doi:10.5061/dryad.4cv1421>
- Lingoes, J. C. 1971. Some boundary conditions for a monotone analysis of symmetric matrices. *Psychometrika*, **36**, 195–203.
- Lloyd, G. T. 2016. Estimating morphological diversity and tempo with discrete character-taxon matrices: implementation, challenges, and future directions. *Biological Journal of the Linnean Society*, **118**, 131–151.
- Marx, F. G. and Fordyce R. E. 2015. Baleen boom and bust: a synthesis of mysticete phylogeny, diversity and disparity. *Royal Society Open Science*, **2**, 140434. doi :10.1098/rsos.140434.
- Oyston, J. W., Hughes, M., Wagner P. J., Gerber S. and Wills, M. A. 2015. What limits the morphological disparity of clades? *Interface focus*, **5**, 20150042. doi:10.1098/rsfs.2015.0042.
- Pérez, D. E. 2018. Phylogenetic relationships of the family Carditidae (Bivalvia: Archiheterodonta). *Journal of Systematic Palaeontology*. doi: 10.1080/14772019.2018.1532463.
- Prentice, K. C., Ruta, M. and Benton, M. J. 2011. Evolution of morphological disparity in pterosaurs. *Journal of Systematic Palaeontology*, **9**, 337–353.
- R Core Team, 2018. R: A language and environment for statistical computing. *R Foundation for Statistical Computing*, Vienna, Austria. URL: <https://www.R-project.org/>.
- Ruta M., Angielczyk K. D., Fröbisch J. and Benton M. J. 2013a. Decoupling of morphological disparity and taxic diversity during the adaptive radiation of anomodont therapsids. *Proceedings of the Royal Society of London Series B, Biological Sciences*, **280**, 20131071. doi:10.1098/rspb.2013.1071.
- Ruta M., Botha-Brink J., Mitchell S. A. and Benton M. J. 2013b. The radiation of cynodonts and the ground plan of mammalian morphological diversity. *Proceedings of the Royal Society of London Series B, Biological Sciences*, **280**, 20131865. doi:10.1098/rspb.2013.1865.
- Smith, A. J., Rosario, M. V., Eiting, T. P. and Dumont, E. R. 2014. Joined at the hip: Linked characters and the problem of missing data in studies of disparity. *Evolution*, **68**, 2386–2400.
- Thorne, P. M., Ruta, M. and Benton M. J. 2011. Resetting the evolution of marine reptiles at the Triassic-Jurassic boundary. *Proceedings of the National Academy of Sciences of the United States of America*, **108**, 8339–8344.
- Toljagic, O. and Butler, R. J. 2013. Triassic-Jurassic mass extinction as trigger for the Mesozoic radiation of crocodylomorphs. *Biology Letters*, **9**, 20130095.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Wills, M. A. 1998. Crustacean disparity through the Phanerozoic: comparing morphological and stratigraphic data. *Biological Journal of the Linnean Society*, **65**, 455–500.

Wills, M. A. 2001. Morphological disparity: a primer. In Adrain, J. M., Edgecombe G. D. and Lieberman, B. S. (Eds.) *Fossils, phylogeny, and form*. Springer, Boston, MA, USA.

Wright, A. M., Lloyd, G. T. and Hillis, D. M. 2016. Modeling character change heterogeneity in phylogenetic analyses of morphology through the use of priors. *Systematic Biology*, **65**, 602-611.

FIGURE CAPTIONS

Figure 1. Scatter plots of distance to centroid against percentage of missing data in taxa of three selected matrices. The distances are calculated as the Euclidean distance in the morphospace with the corresponding number PCos, and are scaled to unit to facilitate comparisons. Regression lines are shown for each data set. “S” denotes a significant Spearman’s correlation between variables and “NS” a non-significant correlation, with $\alpha = 0.05$. The distances are calculated based on the first 3 PCos and for all the PCos.

Figure 2. Histograms showing the percentage of matrices with significant and negative (black), non-significant (light grey), and significant and positive (dark grey) Spearman’s correlations for the GED and the MORD, with $\alpha = 0.05$. The distances are calculated based on the first 3 PCos and all the PCos.

Figure 3. Variation in the Spearman’s correlation coefficient between the percentage of missing data in taxa and their Euclidean distance to the centroid through increasing percentages of explained variance for three selected data matrices. Filled symbols denote significant correlations, while open symbols indicate non-significant correlations, with $\alpha = 0.05$.

Figure 4. Histograms showing the percentage of replications with significant and negative (black), non-significant (light grey), and significant and positive (dark grey) correlations between the Euclidean distance to the centroid and the percentage of missing data in taxa for three selected matrices, calculated for the GED and the MORD, with $\alpha = 0.05$. Approximately 35% of the entries in the matrices were replaced with missing data, and this procedure was repeated 200 times for each matrix and distance measure. The distance to the centroid was calculated using all the PCos.

Figure 5. Disparity values for the studied disparity measures against the proportion of randomly-distributed missing entries added to morphological matrix of Pérez (2018), calculated from distance

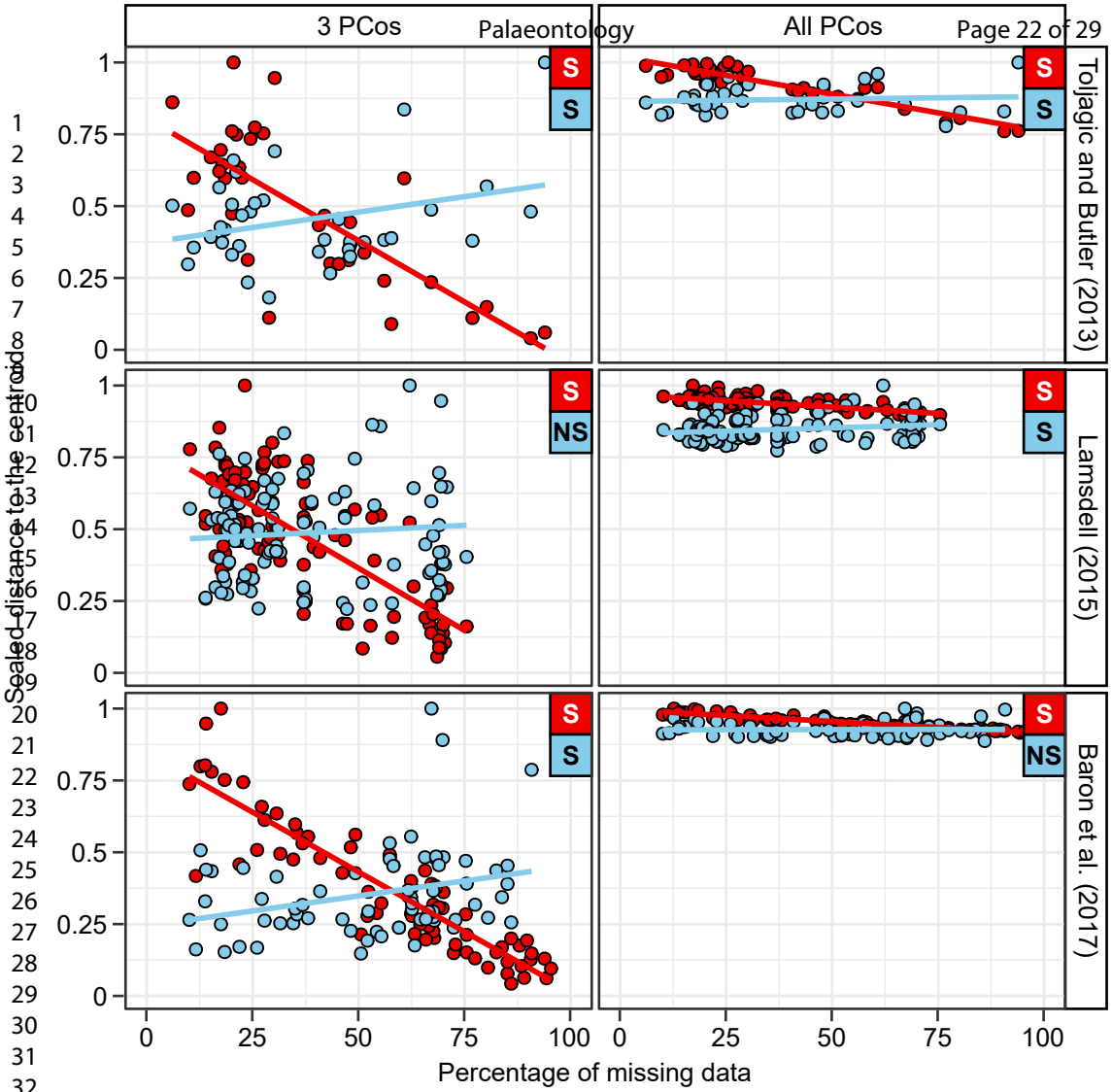
1 matrices generated from the GED and the MORD. Each dot represents the mean of 200 replications.
2
3
4 The values are scaled to unit for each distance and measure to allow an easier visual comparison.
5
6 Vertical bars indicate the 95% confidence interval for the mean. Values of the percentage of missing
7
8 data are treated as different categories, with the GED shifted slightly to the left and the MORD slightly
9
10 to the right in order to avoid overlappings.
11
12
13
14

15 Figure 6. Results of simulations with groups of taxa and different distributions of missing data for the
16
17 matrix of Pérez (2018). Each dot represents the mean of 200 values, with bars indicating the 95%
18
19 bootstrap for the mean. The dashed vertical line represents the reference correlation of 0 (see text for
20
21 a detailed explanation). Values closer to 1 and -1 indicate a stronger influence of the differences in
22
23 the missing data per group in the ranking of the disparity measures. The column ‘%’ indicates the
24
25 percentage of missing data for the matrices in the corresponding simulations and the column ‘S’
26
27 indicates the simulation number according to Table 1.
28
29
30
31
32

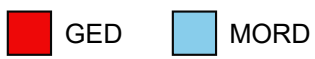
33 Figure 7. Effect of the missing data in a disparity analysis. A: bidimensional morphospace generated
34
35 with the original matrix of Pérez (2018), with only 3.7% of missing or inapplicable entries. B:
36
37 bidimensional morphospace generated after the addition of 35% of missing entries to the taxa marked
38
39 with black dots, producing their displacement (dotted lines) towards the centroid. C: bidimensional
40
41 morphospace generated after the addition of 35% of missing entries to the entire matrix (following a
42
43 beta distribution; same procedure as for the simulations, see text), showing that the spatial structure
44
45 of the groups is diluted. The numbers in parentheses indicate the percentage of variance explained by
46
47 each PCo. The morphospaces of B and C were rotated to match A via Procrustes to accommodate for
48
49 differences in the rotation and scale of the points cloud.
50
51
52
53
54

55 Table 1. Thirteen simulations applied to five randomly-generated groups in selected matrices.
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36

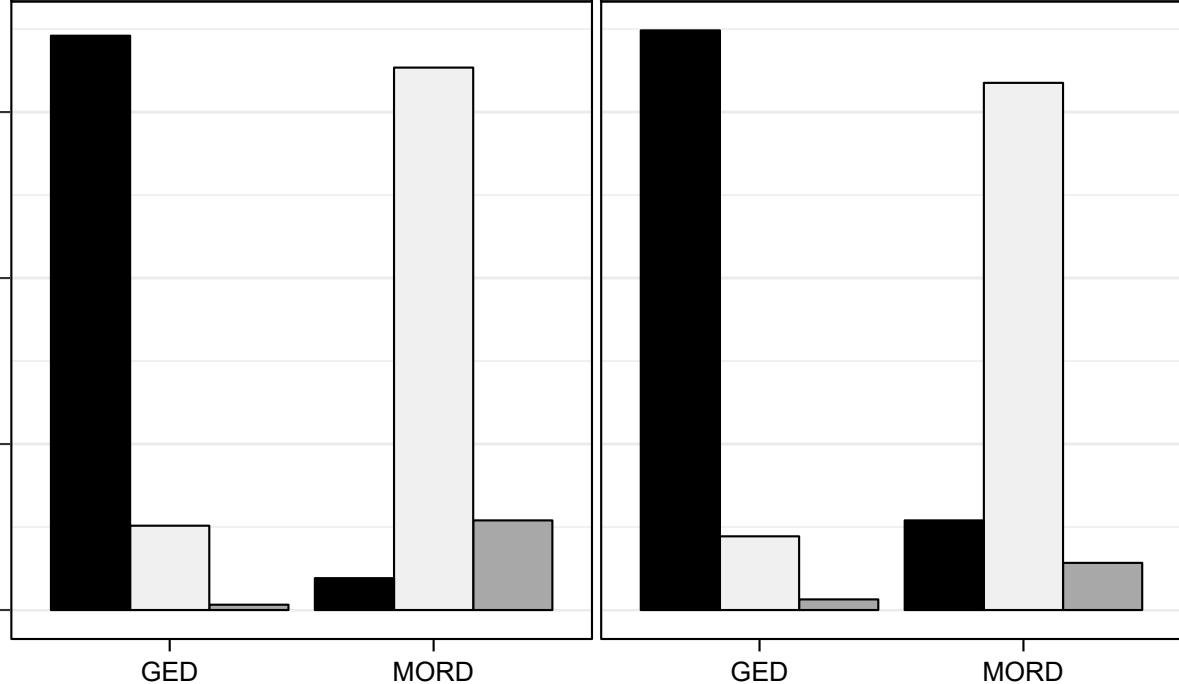


Distance measure Palaeontology



Percentage of matrices

1
2 75
3
4
50
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25



Distance measure

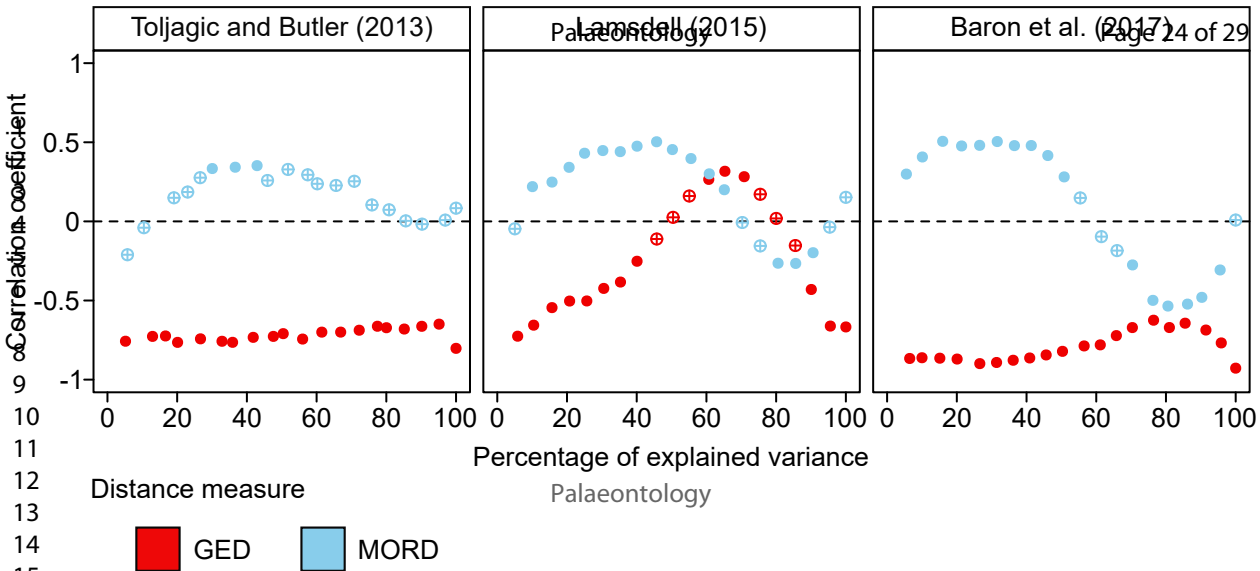
Correlation coefficient

Palaeontology

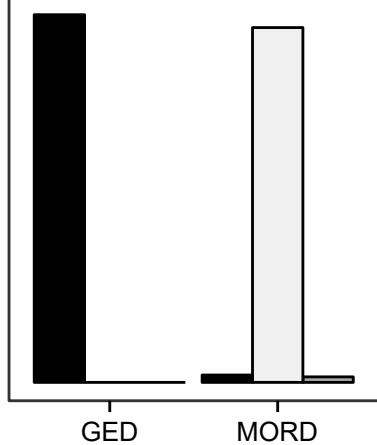
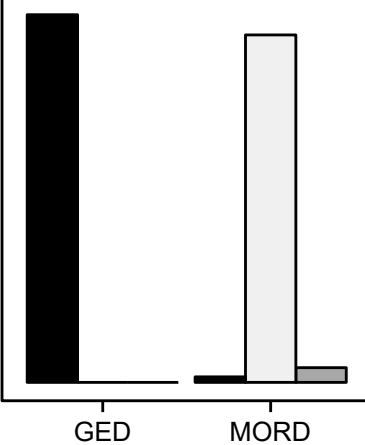
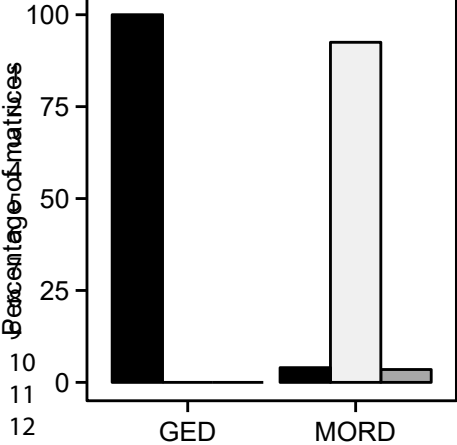
Significant and negative

Not significant

Significant and positive



Percentage of matrices



Distance measure

Correlation coefficient



Significant and negative

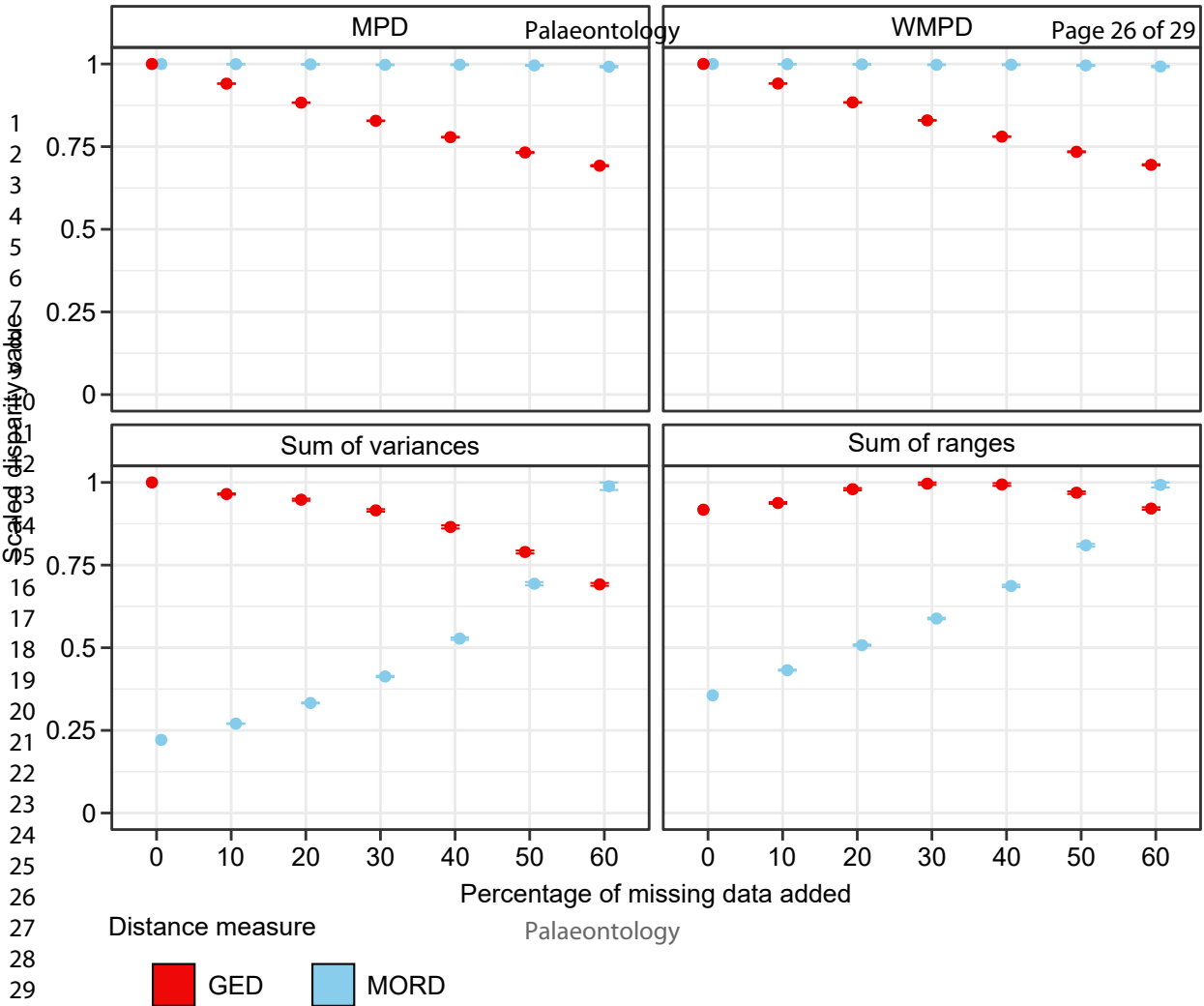


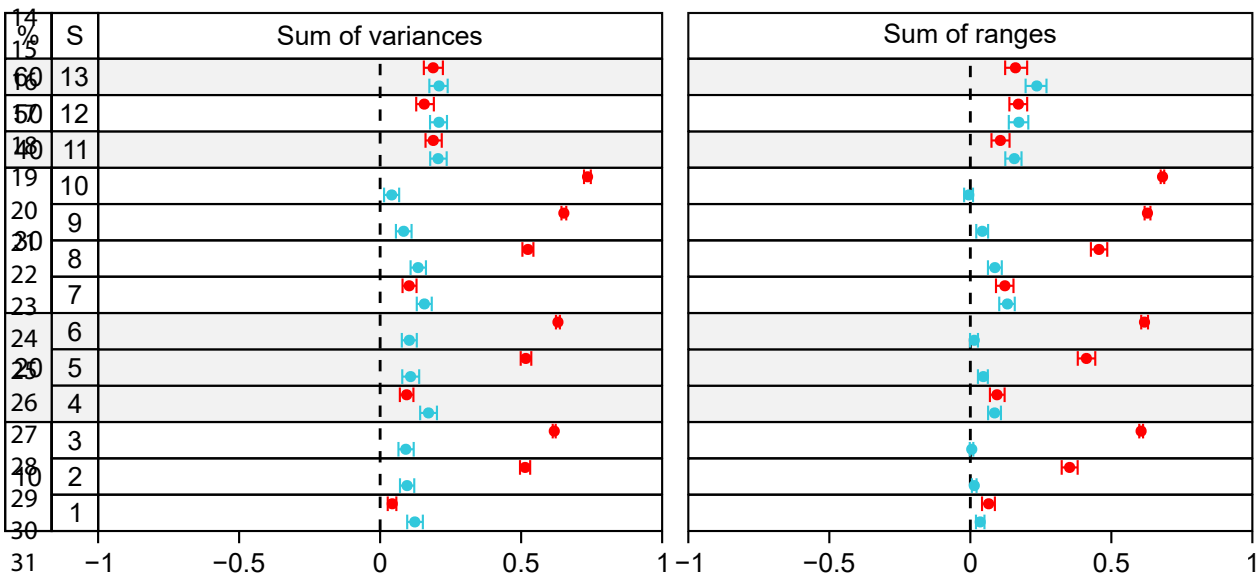
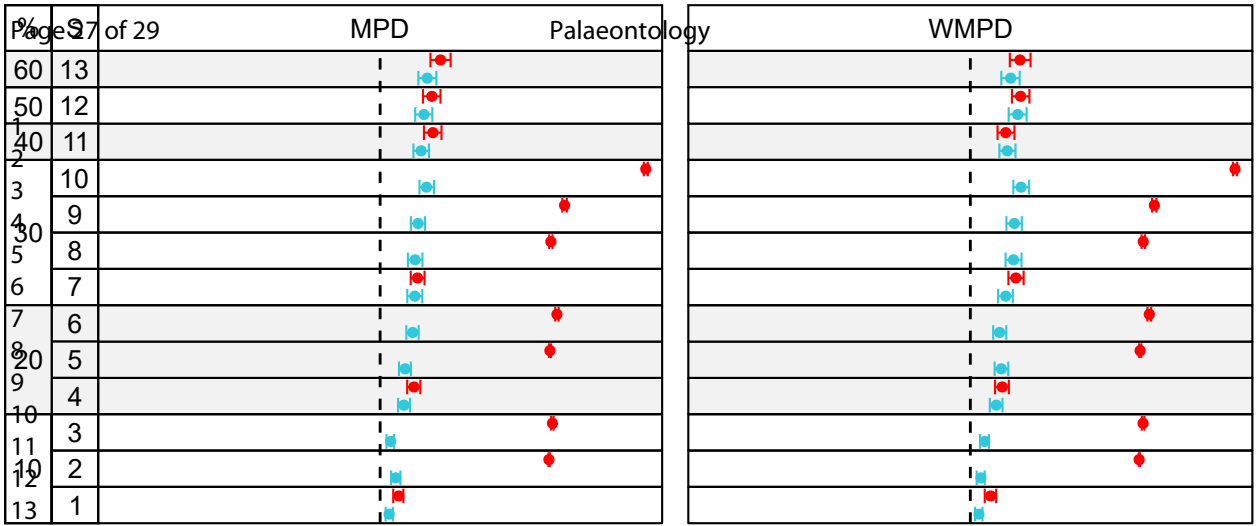
Not significant



Significant and positive

17





14 S

15

16

17

18

19

20

20

22

23

24

25

26

27

28

29

30

31 -1 -0.5 0 0.5 1 -1 -0.5 0 0.5 1

32

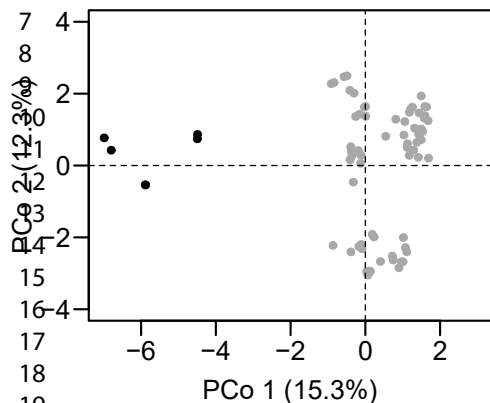
33

34

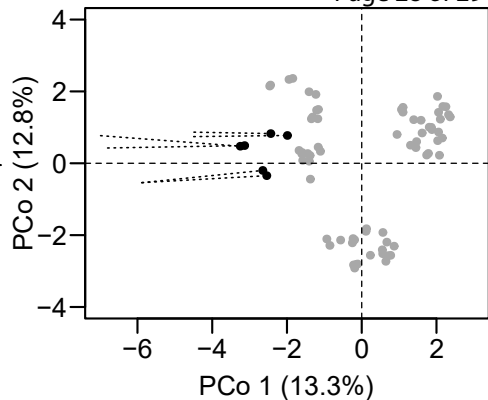
35

36

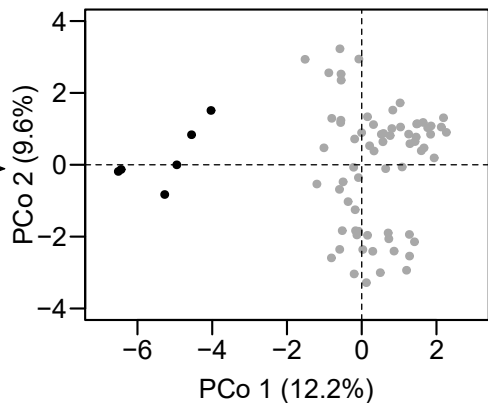
A. Complete matrix



Addition of 35% missing entries
to the black clade



C. Incomplete matrix



Addition of 35% missing entries
to the matrix

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

	Group 1	Group 2	Group 3	Group 4	Group 5	Mean	Step
Simulation 1	10.0%	10.0%	10.0%	10.0%	10.0%	10.0%	0.0%
Simulation 2	15.0%	12.5%	10.0%	7.5%	5.0%	10.0%	2.5%
Simulation 3	20.0%	15.0%	10.0%	5.0%	0.0%	10.0%	5.0%
Simulation 4	20.0%	20.0%	20.0%	20.0%	20.0%	20.0%	0.0%
Simulation 5	25.0%	22.5%	20.0%	17.5%	15.0%	20.0%	2.5%
Simulation 6	30.0%	25.0%	20.0%	15.0%	10.0%	20.0%	5.0%
Simulation 7	30.0%	30.0%	30.0%	30.0%	30.0%	30.0%	0.0%
Simulation 8	35.0%	32.5%	30.0%	27.5%	25.0%	30.0%	2.5%
Simulation 9	40.0%	35.0%	30.0%	25.0%	20.0%	30.0%	5.0%
Simulation 10	50.0%	40.0%	30.0%	20.0%	10.0%	30.0%	10.0%
Simulation 11	40.0%	40.0%	40.0%	40.0%	40.0%	40.0%	0.0%
Simulation 12	50.0%	50.0%	50.0%	50.0%	50.0%	50.0%	0.0%
Simulation 13	60.0%	60.0%	60.0%	60.0%	60.0%	60.0%	0.0%