

# Redescription mining for analyzing local limiting conditions: A case study on the biogeography of large mammals in China and southern Asia

## *Appendix*

Original article published in Ecological Informatics (2021)  
<https://doi.org/10.1016/j.ecoinf.2021.101314>

### AUTHORS:

**Esther Galbrun** (corresponding author)

[esther.galbrun@uef.fi](mailto:esther.galbrun@uef.fi)

School of Computing, University of Eastern Finland,  
Technopolis, Microkatu 1, FI-70210 Kuopio, Finland.

**Hui Tang**

Department of Geosciences, University of Oslo,  
P.O. Box 1022, University of Oslo, NO-0315 Oslo, Norway.

**Anu Kaakinen**

Department of Geosciences and Geography, University of Helsinki,  
P.O. Box 64, FI-00014 University of Helsinki, Finland.

**Indrė Žliobaitė**

Department of Computer Science, University of Helsinki,  
P.O. Box 64, FI-00014 University of Helsinki, Finland.

### SOURCE CODE AND DATA AVAILABILITY:

For this analysis, we only use freely available software and libraries. The datasets used in this study along with the scripts for performing the analysis with classical methods as well as with redescription mining, are publicly available at <https://github.com/zliobaite/redescription-China>.

# A Materials for the case study

## A.1 Delimitation of the study region

The coverage of our dataset is shown in Fig. 1. Localities depicted in green belong to our study area covering China and southern Asia, while other localities included in the global dataset of (Galbrun et al., 2018) are depicted in grey. The study region (green) corresponds to the union of the five rectangles defined by the corners given in Table 1. The rectangles were selected because they demarcate the region we are interested in while also taking advantage of the natural boundaries and the boundaries generated by the requirement on the minimum number of distinct taxa.

Table 1: Coordinates of the corners of the five rectangles used to define the study region.

North-West	20°N, 80°E	10°N, 90°E	5°N, 66°E	28° N, 67°30'E	35°N, 80°E
South-East	35°N, 125°E	20°N, 115°E	28°N, 90°E	37°30'N, 90° E	40N, 120°E

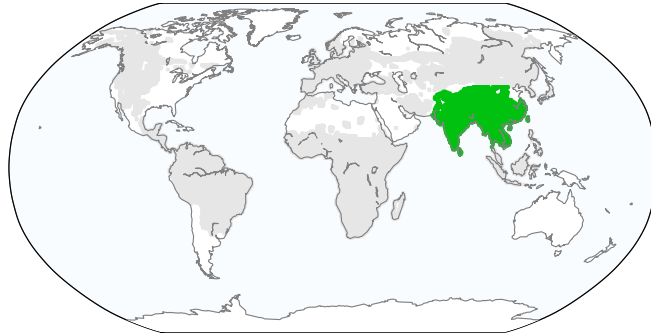


Fig. 1: Maps of the localities contained in the dataset. Localities depicted in green belong to our study area covering China and southern Asia, while other localities included in the global dataset of (Galbrun et al., 2018) are depicted in grey.

## A.2 Datasets

Fig. 2 shows the data aggregation process. Fig. 3, 4 and 5 show the distribution of dental and climatic variables within the study region.

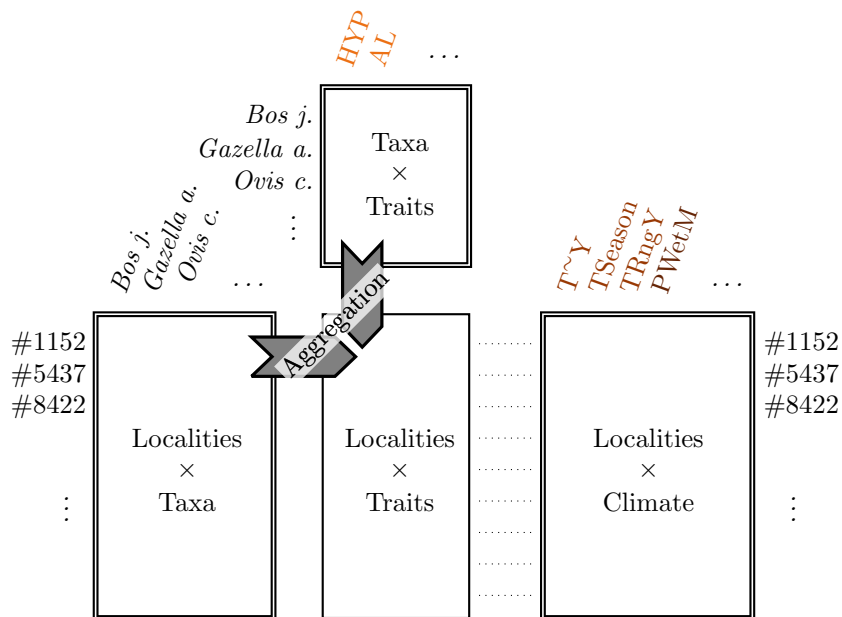


Fig. 2: Datasets and data aggregation. The initial datasets (Localities × Taxa) and (Taxa × Traits) are aggregated to produce the (Localities × Traits) dataset.

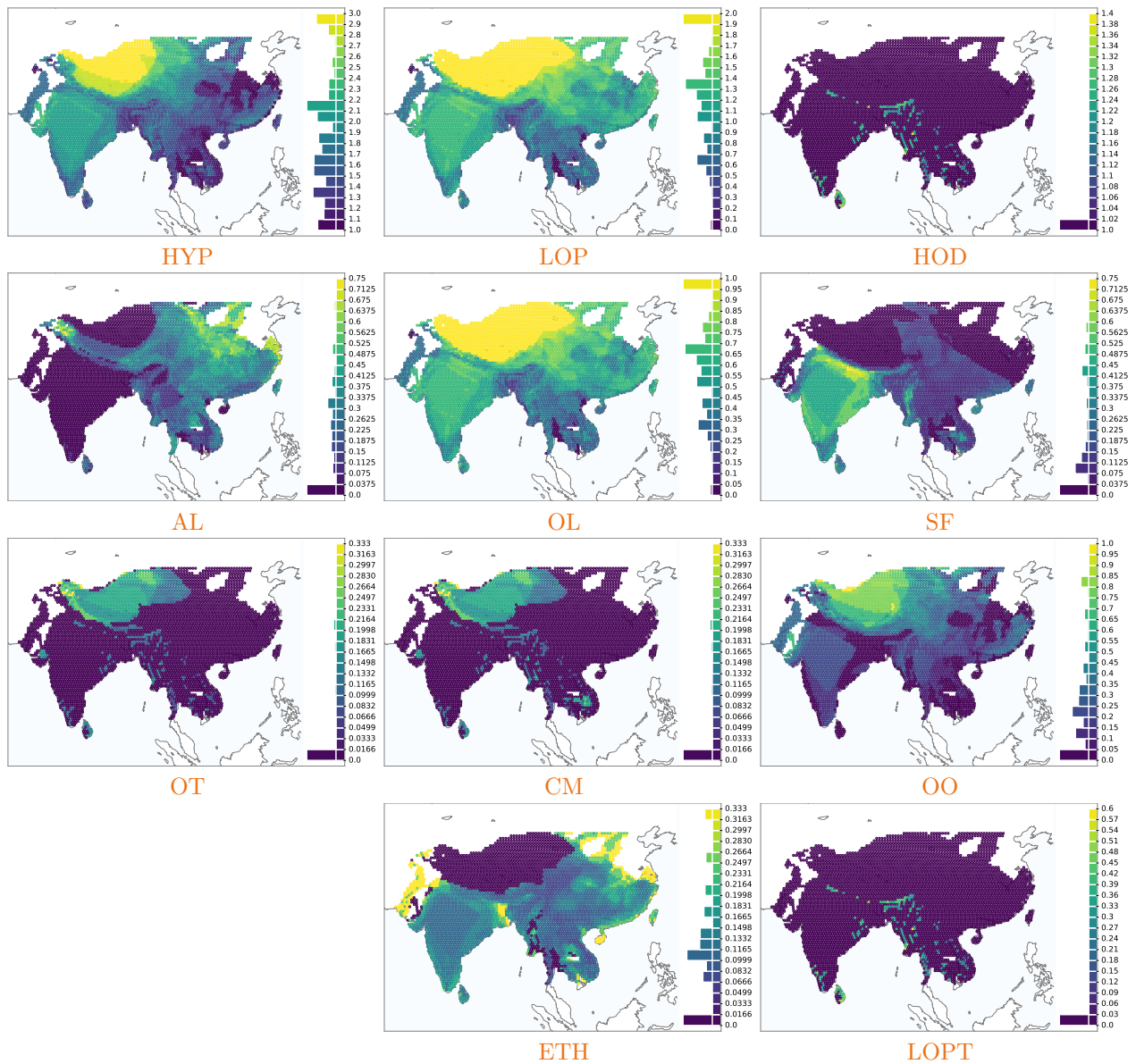


Fig. 3: Maps of dental trait distributions. The color of each dot denotes the value of the trait variable at the corresponding locality. The colored horizontal histogram on the right provides the legends for the colors and indicates the frequency of the different values. The plots are at high resolution intended for electronic viewing and can be zoomed in.

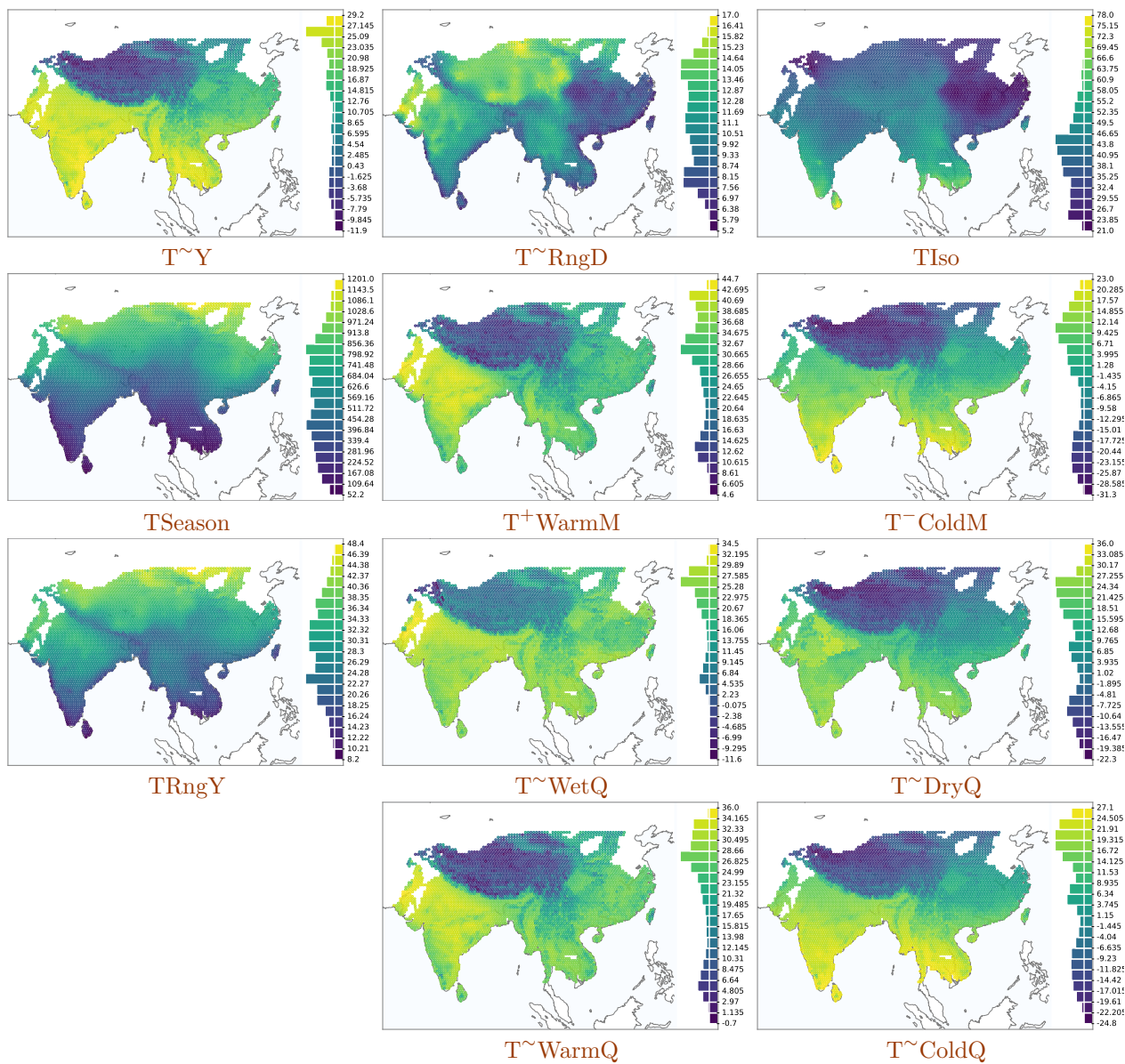


Fig. 4: Maps of bioclimatic variables: temperatures. The color of each dot denotes the value of the temperature variable at the corresponding locality. The colored horizontal histogram on the right provides the legends for the colors and indicates the frequency of the different values. The plots are at high resolution intended for electronic viewing and can be zoomed in.



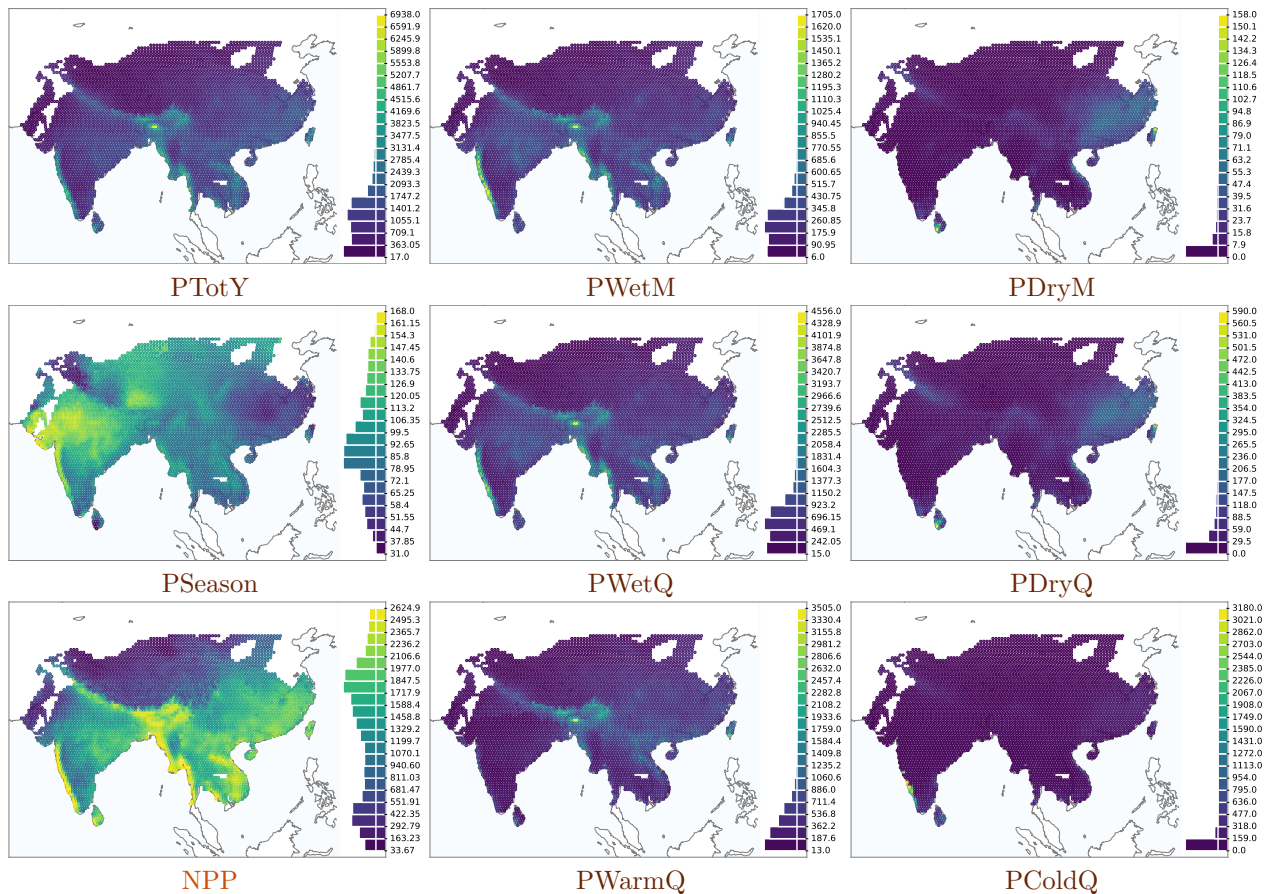


Fig. 5: Maps of bioclimatic variables: precipitation and net primary product (**NPP**, computed from mean annual temperature and precipitation). The color of each dot denotes the value of the variable at the corresponding locality. The colored horizontal histogram on the right provides the legends for the colors and indicates the frequency of the different values. The plots are at high resolution intended for electronic viewing and can be zoomed in.

## B Preamble: classical analysis methods

### B.1 Pairwise correlation and scatter plots

Pearson correlation coefficient and scatter plots for a subset of climate variables are shown in Fig. 6.

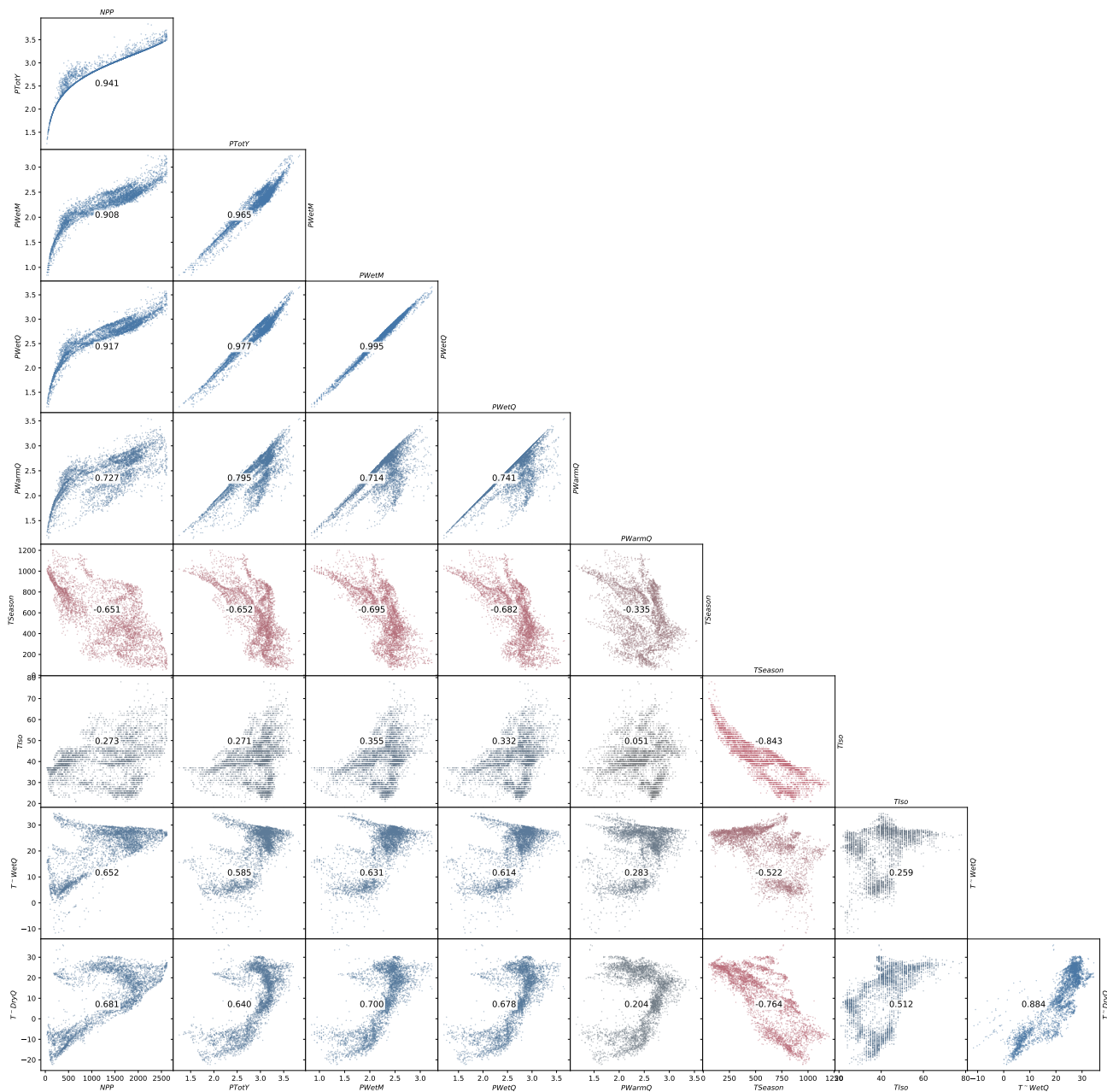


Fig. 6: Pearson correlation coefficient ( $r$ ) and scatter plots for a subset of climate variables. In each plot, the dots represent the localities and are drawn in color according to the Pearson correlation coefficient of the corresponding pair of variables, from blue for strong positive correlation ( $r$  close to 1) to red for strong negative correlation ( $r$  close to  $-1$ ).

### B.2 Multivariate projections

PCA projections considering *Dental traits* and *Climate* variables separately and together are shown in in Fig. 7.



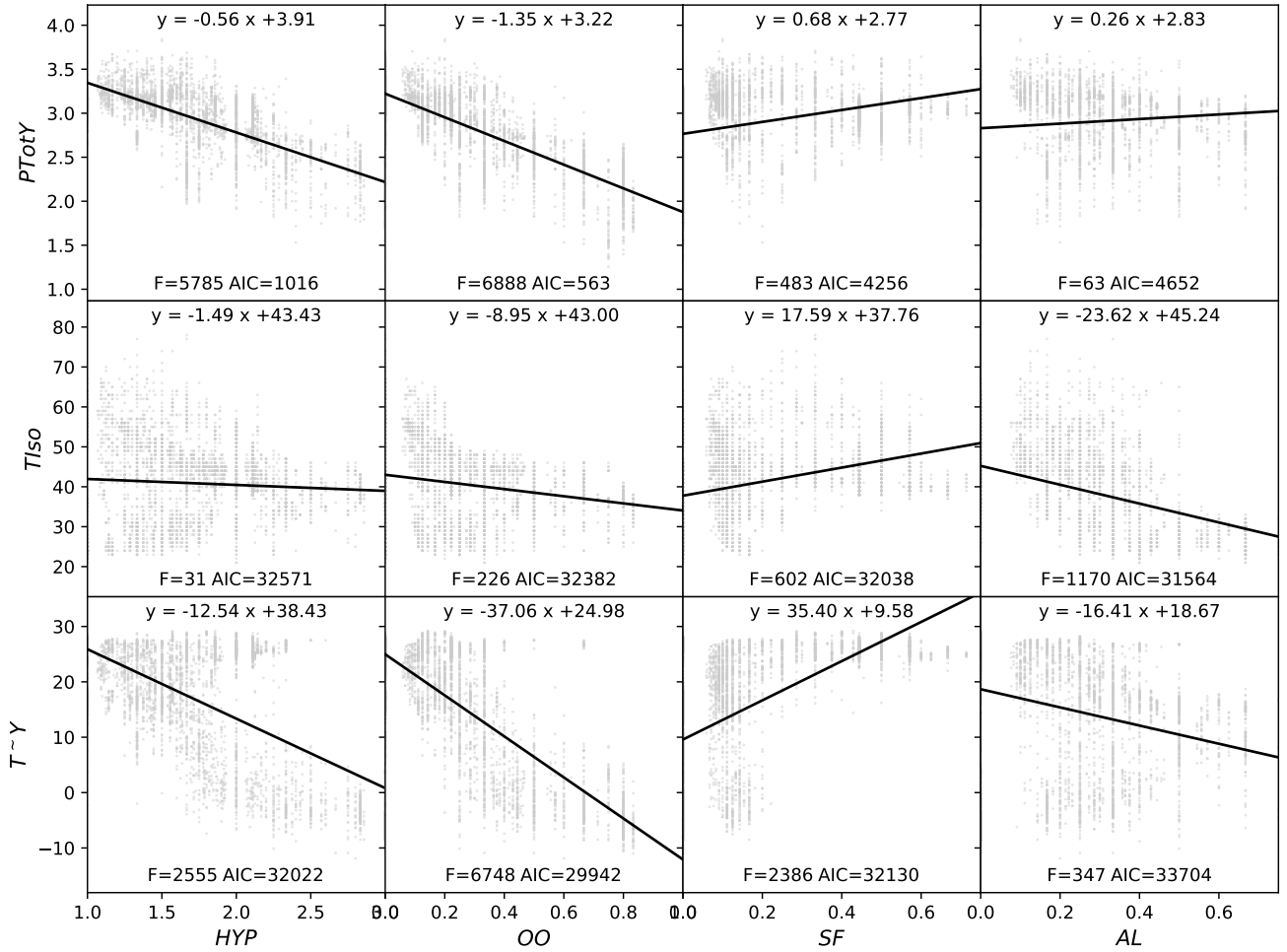


Fig. 8: Regression plots. Estimating one climate variable ( $y$ ) as a linear function of a dental trait variable ( $x$ ) using the generalized least squares (GLS) (Aitken, 1934). For each model, we indicate the corresponding equation (top) as well Akaike's information criterion (AIC) and the F-statistic (F) of the model (bottom).

### B.3 Regression models

The best-fitting linear regression models according to Akaike's information criterion (AIC) for  $PTotY$ ,  $TIso$  and  $T\sim Y$  are listed as examples in Table 2. Examples of single-variate models are plotted in Fig. 8.

### B.4 Clustering

The distance  $D(U, V)$  between two clusters  $U$  and  $V$  is defined as follows in different HCA methods, also referred to as *linkage functions*, including:

- single** (SL), the minimum distance between cluster members, i.e.  $D(U, V) = \min_{(u,v) \in U \times V} d(u, v)$
- complete** (CL), the maximum distance between cluster members, i.e.  $D(U, V) = \max_{(u,v) \in U \times V} d(u, v)$
- average** a.k.a. Unweighted Pair-Group Method using arithmetic Averages (UPGMA), the average distance between cluster members, i.e.  $D(U, V) = \sum_{(u,v) \in U \times V} d(u, v) / (|U| \cdot |V|)$
- weighted** a.k.a. Weighted Pair-Group Method using arithmetic Averages (WPGMA), a variant of UPGMA weighted by the size of the clusters
- centroid** a.k.a. Unweighted Pair-Group Method using Centroids (UPGMC), the distance between cluster centroids, i.e.  $D(U, V) = d(c_U, c_V)$
- median** a.k.a. Weighted Pair-Group Method using Centroids (WPGMC), a variant of UPGMC weighted by the size of the clusters
- ward** relies on Ward's minimum variance criterion

The *silhouette coefficient* was introduced by (Rousseeuw, 1987) to evaluate the quality of a clustering. We consider a set of data points  $E$ , in our case localities, and a clustering  $\{C_1, C_2, \dots\}$  such that each data point



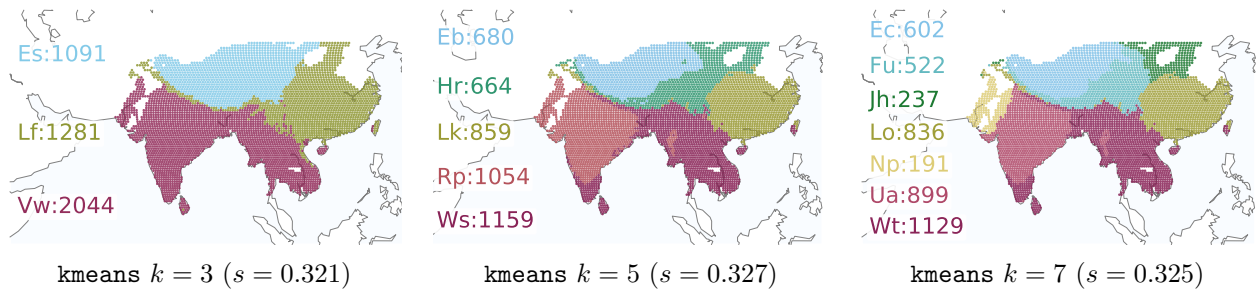


Fig. 9:  $k$ -means clusterings and associated silhouette coefficients ( $s$ ). To the left of each map, we list the different clusters, with the number of localities they contain. Labels and colors are assigned to clusters in such a way that more similar clusters (across all generated clusterings) receive similar colors and labels that are close in the alphabetical order.

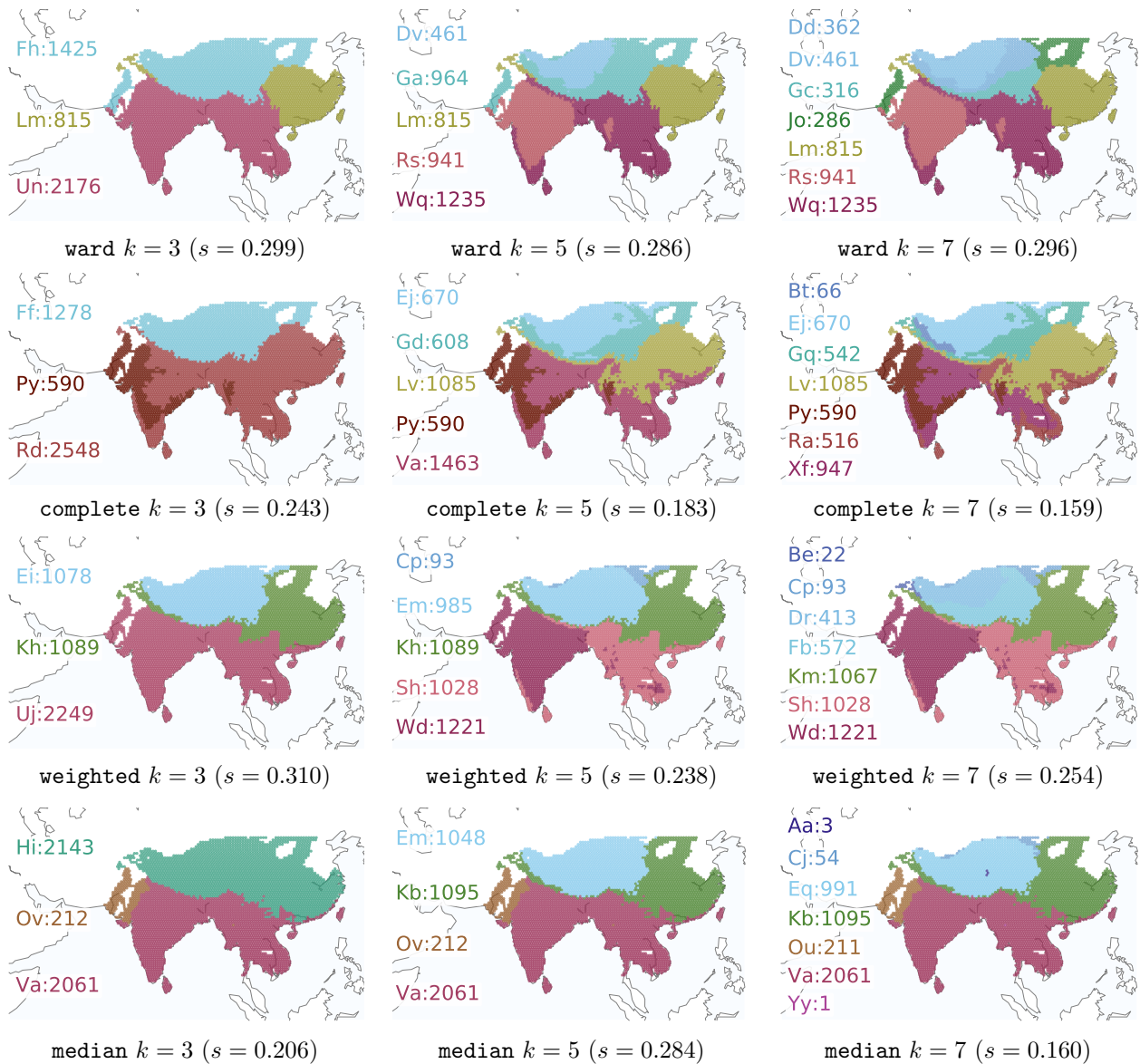


Fig. 10: Hierarchical agglomerative clusterings and associated silhouette coefficients ( $s$ ). Each panel corresponds to a particular combination of linkage function and number of clusters ( $k$ ) as indicated underneath. To the left of each map, we list the different clusters, with the number of localities they contain. Labels and colors are assigned to clusters in such a way that more similar clusters (across all generated clusterings) receive similar colors and labels that are close in the alphabetical order.



belongs to one cluster. Furthermore, we consider a function  $d$  that measures the distance between any pair of data points with respect to data variables.

For data point  $x$  belonging to cluster  $C_i$ , let

$$a(x) = \frac{1}{|C_i| - 1} \sum_{y \in C_i, x \neq y} d(x, y) \quad \text{and} \quad b(x) = \min_{k \neq i} \frac{1}{|C_k|} \sum_{j \in C_k} d(i, j).$$

That is,  $a(x)$  denotes the average distance between  $x$  and other points in the same cluster whereas  $b(x)$  is the minimum, over the other clusters, of average distance between  $x$  and points in that cluster.

Then, the silhouette value for point  $x \in C_i$  is

$$s(x) = \begin{cases} \frac{b(x) - a(x)}{\max(a(x), b(x))} & \text{if } |C_i| > 1, \\ 0 & \text{otherwise.} \end{cases}$$

and the silhouette coefficient is the average silhouette value over all data points

$$s = \sum_{x \in E} \frac{s(x)}{|E|}.$$

Clusterings obtained with the  $k$ -means algorithm and with hierarchical agglomerative algorithms for  $k = 3, 5$  and 7 clusters can be found in Fig. 9 and 10, respectively.

Note that the merging process in hierarchical agglomerative algorithms may result in fewer clusters than specified, as is the case here with the **median** variant where four clusters are returned when setting  $k = 5$ . Single linkage as well as average linkage and centroid linkage (i.e. unweighted variants) returned partitions that typically consisted of one very large cluster together with extremely small clusters (e.g. a dozen localities each), which is not very informative for further analysis. Obtaining a well-balanced hierarchical clustering is often methodologically challenging, and here it may be further complicated by the presence of spatial autocorrelation in the data.

To facilitate the visual comparison of the (many) clusters produced by the different methods, we compute a one-dimensional distance-preserving projection of the clusters (by building the graph of pair-wise similarities between the clusters and computing the vector associated to the second-smallest eigenvalue of the associated Laplacian matrix, a.k.a. Fiedler vector, commonly used to partition graphs) and assign colors along a gradient, in such a way that similar clusters are close to each other and assigned similar colors.

## C Redescription mining methodology

To assess the statistical significance of a redescription, we compute a  $p$ -value that indicates how likely it is that the support of the redescription is as large or larger than observed, given the size of the support of the two queries it consists of, assuming the queries are independent. Consider two statistically independent random queries with marginal probabilities equal to those of the queries under consideration, i.e. with marginal probabilities equal to the fraction of covered localities  $P(q) = |\text{supp}(q)| / |E|$ . The  $p$ -value is computed using the binomial distribution, as

$$\text{pval}(q_{\mathbf{D}}, q_{\mathbf{C}}) = \sum_{s=|S|}^{|E|} \binom{|E|}{s} (P(q_{\mathbf{D}})P(q_{\mathbf{C}}))^s (1 - P(q_{\mathbf{D}})P(q_{\mathbf{C}}))^{|E|-s},$$

where  $S = \text{supp}((q_{\mathbf{D}}, q_{\mathbf{C}}))$  is the observed support of the redescription. The  $p$ -value is the probability of obtaining a set of same size or larger if each element of a set of size  $|E|$  has a probability equal to the product of marginals  $P(q_{\mathbf{D}})$  and  $P(q_{\mathbf{C}})$  to be selected, in accordance with the independence assumption.

## D Case study: biogeographic analysis with redescription mining

For reference, the top ten redescriptions from the two runs are listed in Table 3.

Summaries obtained with for the different clustering variants and for  $k = 3, 5$  and 7 clusters are shown in Fig. 11. Since they are based on the support of redescriptions, the clusterings naturally take into account both dental traits and bioclimatic variables, whereas when clustering the raw data, the different nature and scale of the variables might raise issues.

Table 3: Top-ten redescrptions from the two runs. For each redescription, we list its queries, that is, the query over dental traits variables ( $q_D$ ) and the query over bioclimatic variables ( $q_C$ ). We also indicate the accuracy of the redescription (J) as well as the size of its support, as the number of localities described ( $|\text{supp}|$ ) and as a percentage of the total number of localities ( $\text{supp}\%$ ).

<b>R1.1</b>	J = 0.846  supp  = 2540 supp % = 57.52 $q_D = [\text{LOP} \leq 1.556] \text{ AND } [\text{OO} \leq 0.25]$ $q_C = [19.9 \leq \text{T}^{\sim}\text{WetQ} \leq 29.6] \text{ AND } [388 \leq \text{PTotY}]$	<b>R2.1</b>	J = 0.836  supp  = 741 supp % = 16.78 $q_D = [1.8 \leq \text{LOP}] \text{ AND } [\text{ETH} \leq 0]$ $q_C = [\text{T}^{\sim}\text{WarmQ} \leq 17.9] \text{ AND } [\text{PTotY} \leq 496]$
<b>R1.2</b>	J = 0.845  supp  = 893 supp % = 20.22 $q_D = [0.769 \leq \text{OL}] \text{ AND } [\text{SF} \leq 0.222] \text{ AND } [\text{ETH} \leq 0.125]$ $q_C = [\text{T}^{\sim}\text{Y} \leq 5.9] \text{ AND } [592.8 \leq \text{TSeason} \leq 1064]$	<b>R2.2</b>	J = 0.829  supp  = 2553 supp % = 57.81 $q_D = [\text{OO} \leq 0.25]$ $q_C = [19.9 \leq \text{T}^{\sim}\text{WetQ} \leq 29.6]$
<b>R1.3</b>	J = 0.810  supp  = 2292 supp % = 51.90 $q_D = [0.222 \leq \text{LOP} \leq 1.333] \text{ AND } [\text{AL} \leq 0.417]$ $q_C = [7.6 \leq \text{T}^{\sim}\text{DryQ}] \text{ AND } [6.5 \leq \text{T}^{\sim}\text{ColdQ} \leq 25.2]$	<b>R2.3</b>	J = 0.822  supp  = 1863 supp % = 42.19 $q_D = [\text{LOP} \leq 0.769] \text{ OR } [0.167 \leq \text{SF}]$ $q_C = [16.8 \leq \text{T}^{\sim}\text{DryQ}]$
<b>R1.4</b>	J = 0.810  supp  = 2271 supp % = 51.43 $q_D = [\text{HYP} \leq 1.889] \text{ AND } [\text{LOP} \leq 1.75] \text{ AND } [\text{OO} \leq 0.4]$ $q_C = [19.6 \leq \text{T}^+\text{WarmM} \leq 38.5] \text{ AND } [116 \leq \text{PWarmQ}]$	<b>R2.4</b>	J = 0.808  supp  = 2494 supp % = 56.48 $q_D = [0.222 \leq \text{LOP} \leq 1.333]$ $q_C = [15 \leq \text{T}^{\sim}\text{Y} \leq 27.4]$
<b>R1.5</b>	J = 0.787  supp  = 1982 supp % = 44.88 $q_D = [0.222 \leq \text{LOP} \leq 1.833] \text{ AND } [0.077 \leq \text{AL} \leq 0.667]$ $\text{AND } [\text{OT} \leq 0.111]$ $q_C = [3.6 \leq \text{T}^{\sim}\text{WarmQ} \leq 29.4]$ $\text{AND } [266 \leq \text{PWarmQ} \leq 1556]$	<b>R2.5</b>	J = 0.783  supp  = 2315 supp % = 52.42 $q_D = [\text{HYP} \leq 1.889]$ $q_C = [19.6 \leq \text{T}^+\text{WarmM} \leq 38.5]$
<b>R1.6</b>	J = 0.769  supp  = 2025 supp % = 45.86 $q_D = [\text{HYP} \leq 2.125] \text{ AND } [\text{LOP} \leq 1.5]$ $\text{AND } [0.059 \leq \text{SF} \leq 0.571]$ $q_C = [\text{TRngY} \leq 30.3] \text{ AND } [-1.8 \leq \text{T}^{\sim}\text{ColdQ} \leq 26.6]$	<b>R2.6</b>	J = 0.749  supp  = 2306 supp % = 52.22 $q_D = [0.077 \leq \text{AL}]$ $q_C = [6.2 \leq \text{T}^{\sim}\text{WarmQ} \leq 28.5]$
<b>R1.7</b>	J = 0.751  supp  = 1959 supp % = 44.36 $q_D = [\text{AL} \leq 0.417] \text{ AND } [0.059 \leq \text{SF}]$ $\text{AND } [0.059 \leq \text{ETH} \leq 0.167]$ $q_C = [92.7 \leq \text{TSeason} \leq 646]$ $\text{AND } [631.57 \leq \text{NPP} \leq 2608.99]$	<b>R2.7</b>	J = 0.717  supp  = 2212 supp % = 50.09 $q_D = [0.111 \leq \text{ETH}]$ $q_C = [29.1 \leq \text{T}^+\text{WarmM}]$
<b>R1.8</b>	J = 0.749  supp  = 2210 supp % = 50.05 $q_D = [0.077 \leq \text{AL}] \text{ AND } [\text{SF} \leq 0.25]$ $q_C = [-5.5 \leq \text{T}^{\sim}\text{Y} \leq 23.9] \text{ AND } [\text{T}^+\text{WarmM} \leq 34.8]$	<b>R2.8</b>	J = 0.708  supp  = 2185 supp % = 49.48 $q_D = [0.083 \leq \text{ETH} \leq 0.2]$ $q_C = [181 \leq \text{PWetM}]$
<b>R1.9</b>	J = 0.735  supp  = 2168 supp % = 49.09 $q_D = [\text{LOP} \leq 1.556] \text{ AND } [\text{OO} \leq 0.4]$ $\text{AND } [0.083 \leq \text{ETH} \leq 0.2]$ $q_C = [\text{TSeason} \leq 870.3] \text{ AND } [181 \leq \text{PWetM}]$	<b>R2.9</b>	J = 0.685  supp  = 2109 supp % = 47.76 $q_D = [0.059 \leq \text{SF} \leq 0.571]$ $q_C = [\text{TRngY} \leq 30.3]$
<b>R1.10</b>	J = 0.732  supp  = 2045 supp % = 46.31 $q_D = [\text{HYP} \leq 2.444] \text{ AND } [0.077 \leq \text{AL}] \text{ AND } [\text{SF} \leq 0.375]$ $q_C = [14.3 \leq \text{T}^+\text{WarmM} \leq 34.6] \text{ AND } [1 \leq \text{PDryM}]$	<b>R2.10</b>	J = 0.681  supp  = 2039 supp % = 46.17 $q_D = [\text{LOP} \leq 1.286]$ $q_C = [\text{T}^{\sim}\text{RngD} \leq 12.3]$

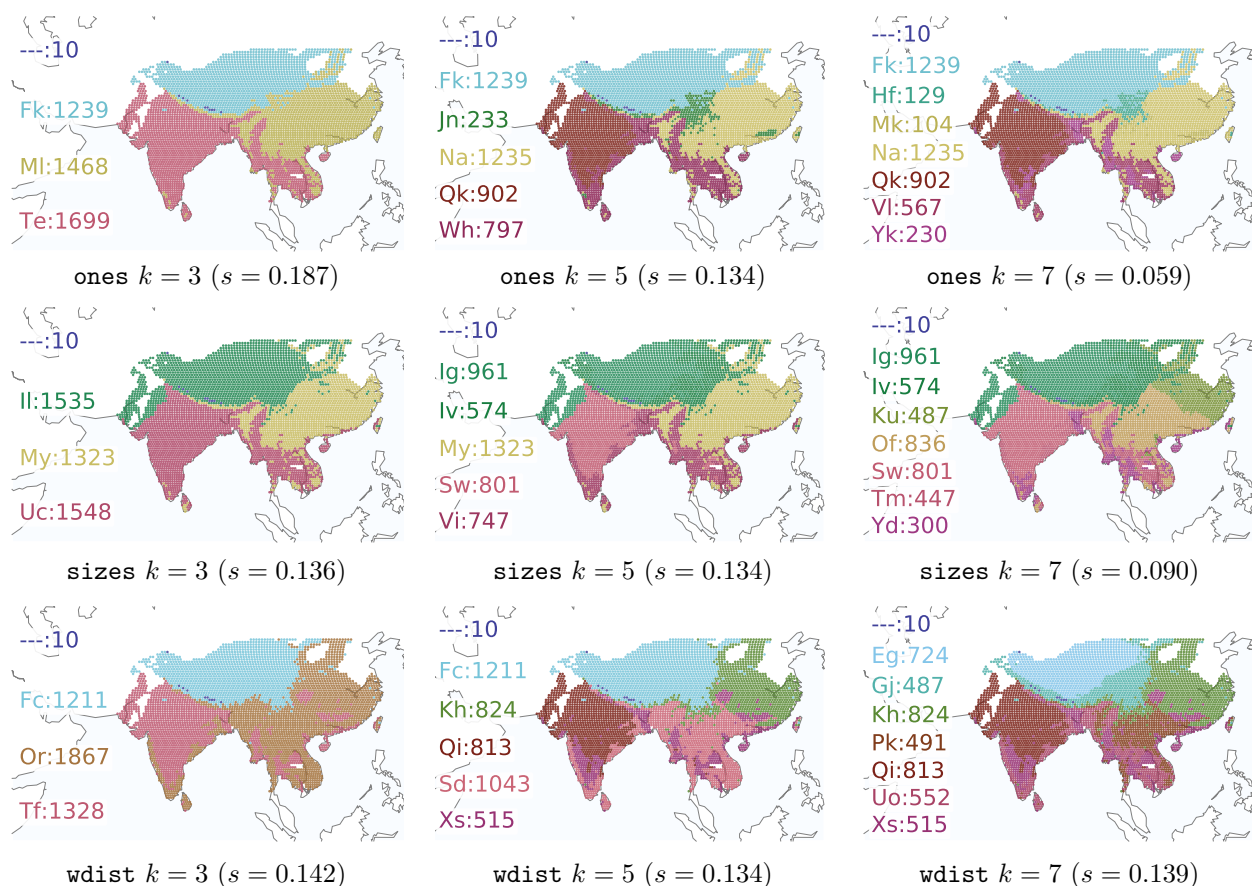


Fig. 11: Redescription summaries, i.e. redescription-based clusterings, and associated silhouette coefficients ( $s$ ). To the left of each map, we list the different clusters, with the count of localities they contain. Labels and colors are assigned to clusters in such a way that more similar clusters (across all generated clusterings) receive similar colors and labels that are close in the alphabetical order. There are a few localities (in dark blue) that do not support any of the selected redescriptions and hence do not belong to any cluster.

Compared to the results of the clustering analysis based on raw dental traits and climate data (cf. Fig. 9 and 10), the redescription summaries (cf. Fig. 11) reveal generally similar distinct biogeographical regions, such as the Tibetan Plateau, East China and India. However, they exhibit a lower similarity between India and Southeast Asia and between the Tibetan Plateau and northern China, but a greater similarity between southern China and Southeast Asia and between southern China and northern China, especially when considering few clusters (i.e. for  $k = 3$  or  $5$ ). With more clusters ( $k = 5$  or  $7$ ), the results based on redescriptions reveal a much finer spatial structure over southern China and Southeast Asia, which corresponds well to plant relicts found in these regions Huang et al. (2015). In contrast, the results based on raw data seem to focalize on the differences within the Tibetan Plateau. These discrepancies exemplify the potential added value of using redescriptions-based clusters to delineate biologically meaningful ecoregions.

## References

- Aitken, A. C. (1934). “On least squares and linear combination of observations”. In: *Proceedings of the Royal Society of Edinburgh*. Section B: Biological Sciences 55, pp. 42–48.
- Galbrun, E., H. Tang, M. Fortelius, and I. Žliobaitė (2018). “Computational biomes: The ecometrics of large mammal teeth”. In: *Paleontologia Electronica*. DOI: 10.26879/786.
- Huang, Y., F. Jacques, T. Su, D. Ferguson, H. Tang, W. Chen, and Z. Zhou (2015). “Distribution of Cenozoic plant relicts in China explained by drought in dry season”. In: *Scientific Reports* 5, p. 14212. DOI: 10.1038/srep14212.
- Rousseeuw, P. (1987). “Silhouettes: A graphical aid to the interpretation and validation of cluster analysis”. In: *Journal of Computational and Applied Mathematics* 20, pp. 53–65. DOI: 10.1016/0377-0427(87)90125-7.