

Global ecology in a data-rich world

see CONTRIBUTORS

work in progress

Ecologists are tasked with providing information on issues at the global scale. Notable examples include, in addition to the global biodiversity crisis (???), predicting the consequences of the loss of trophic structure Estes et al. (2011), rapid shifts in species distributions (???), and . Most of these pressing topics require to be addressed (i) at the global scale and (ii) through the integration of several types of data (Thuiller et al. 2013). Because of these requirements, new sampling is not a viable solution: there is no funding structure to finance it, and there are time and scale constraints involved that make it unrealistic.

Thankfully, developments in the recent years means that ecologists can now leverage existing data, and use them to *build* new datasets suitable for the questions at hand. There are several parallel advances that make this approach possible. First, the volume of data on ecological systems that are available *openly* increases on a daily basis. This includes point-occurrence data, as in *e.g.* GBIF, ITIS, but also taxonomic knowledge (through NCBI or EOL), and trait and interactions data. A vast treasure trove of ecological information is now available without having to contact every contributor individually. Second, this data is available in a *programmatic* way. As opposed to manual collection, identification, and curation of datasets, most of these services implement web API that allow to query them, either once or on a regular basis, to retrieve records with the desired properties. This ensures that the process is repeatable, testable, transparent, and error-proof. Finally, most of the heavy-lifting for these tasks can be done through a *burgeoning ecosystem of packages and software*, that take care of handling query formatting, data retrieval, etc, and expose simple interfaces to researchers.

To us, this opens no less than a new area of research for ecologists interested in asking questions at large spatial and organizational scales – we live in a data-rich world, and a very large amount of these data can now *easily* be collected, reducing the need for additional costly and time-consuming sampling. More importantly, this allows *rapid* evaluation of scenarios. In this contribution, building on a real-life example, we (i) outline the basic approach, (ii) identify technical bottlenecks, (iii) discuss ethical considerations, and (iv) provide clear recommendations moving forward.

27 An illustrative case-study

28 Food-web data, that is the determination of trophic interactions among species, are difficult to collect. The usual
29 approach is to assemble literature data, expert knowledge, and additional information coming from field work, either
30 as observation of direct feeding events, or through the analysis of gut content. Because of these technical constraints,
31 food-web data are most often assembled in a single location. This impedes our ability to address the variation of
32 their structure in space, which may both translate the action of macro-ecological mechanisms, as well as hold key to
33 our ability to predict the spatial variation of ecological properties.

34 In this case study, we are interested in predicting the structure of the pine-marsh food web worldwide. One example of
35 this system has been described in detail by (???). We will show how coupling these data with additional interactions
36 reported in the literature, as well as species occurrence data, allows building a predictive model that describes the
37 spatial variation of this system.

38 Due to the limitations, ... (qualitative rather than quantitative prediction)

39 Interactions data

40 Data from the pine-marsh food web were take from **REF**, as made available in the **IWDB** database (URL). Marshes,
41 as almost all wetlands, are critically endangered and home to a host of endemic biodiversity (Fensham et al. 2011,
42 Minckley et al. 2013). They represent a prime example of ecosystems for which data-based prediction can be used to
43 generate scenarios at a temporal scale relevant for conservation decisions, and faster than what sampling could allow.

44 The original food web (105 nodes, including vague denominations like *Unidentified detritus* or *Terrestrial invertebrates*),
45 where cleaned in the following way. First, all nodes where aggregated to the *genus* level. Due to high level of structure
46 in trophic interactions emerging from taxonomic rank alone (Eklof et al. 2011), aggregating to the genus level has
47 the double advantage of (i) removing ambiguities on the identification of species and (ii) allowing to integrate data
48 when any two species from given genera interact. Second, all nodes that where not identified (**Unidentified** or
49 **Unknown** in the original data). The cleaned network documented 227 interactions, between 80 genera. Using the
50 name checking functions from the **taxize** package (Chamberlain and Szöcs 2013) revealed that all of these genus
51 names where valid.

52 Because the original foodweb was sampled *locally*, there is the possibility that interactions between genera are not
53 reported. To circumvent this, we queried the *GLOBI* database (Poelen et al. 2014) for each genus name, and
54 retrieved all *feeding* interactions. For all *new* genera retrieved through this method, we retrieved their interactions
55 with genera already in the network. This network has 789 genera, and a total of 9328 interactions.

56 Occurrence data and filtering

57 For each genera, we downloaded the known occurrences from GBIF and BISON. This yielded 64763 point-presence
58 records. Because the goal is, ultimately, to do spatial modeling of the structure of the network, we removed genera
59 for which less than 100 occurrences were known. This seems like a stringent filter, yet it enables to (i) maintain
60 sufficient predictive powers for SDMs, and (ii) only work on the genera for which we have “high-quality” data.
61 Genera with less than 100 records were removed from the occurrence data and from the metanetwork. The final
62 metanetwork has 4271 interactions between 188 genera. Given the curated publicly available data, it represents the
63 current best description of feeding interactions between species of pine forests. A visual depiction of the network is
64 given in *Fig. 1*.

65 In itself, the fact that filtering for genera with over 100 records reduced the sample size from 739 genera to 188
66 indicates how crucial it is that observations are reported in public databases. The type of analysis we present here is
67 only as good as the underlying data.

68 SDM

69 Assembly

70 Opportunities

71 Hypothesis testing for large-scale systems is inherently limited by the availability of suitable datasets. Perhaps as a
72 result, macro-ecology has been guided by a search for patterns that are very broad both in scale and nature (???), as
73 opposed to the testing of pre-established hypothesis. While it is obvious that collecting data at scales that are large
74 enough to be relevant is an insurmountable effort (both because of the monetary, time, and human costs needed),
75 we suggest that macroecologists should build on existing databases, and aggregate them in a way that allows direct
76 testing of proposals stemming from theory.

77 Challenges

78 The merging of large databases is already asking the question of proper attribution of data paternity. Namely, there
79 are two core issues that need community consultation in order to be resolved. First, *what is the proper mode of*
80 *attribution when a very large volume of data is aggregated?* Second, *what should be the intellectual property of the*
81 *newly aggregated dataset?*

82 This approach hardly qualifies as *big data*; nevertheless, it relies on the management and integration of a large
83 volume of heterogeneous information, both qualitatively larger than the current “norm”. The first challenge is
84 being able to *manage* this data; it requires data management skills that are not usually needed when ... **REF**.
85 The second challenge is being able to *manipulate* these data; small though they may be, the data from this case
86 study do not fit in the memory of R (arguably the most commonly known software in ecology) without some
87 adjustments. It is now worth asking whether our total reliance on this tool (as opposed to more performing yet
88 equally user-friendly languages as `python` or `julia`) is going to pay-off in the long term. Since it makes very little
89 doubt that computational approaches are going to become increasingly common in ecology **REF**, and are one areas
90 where the community identifies lower skills that would be needed (Barraquand et al. 2014), it seems evident that
91 developing computational literacy should be a part of the core ecological curriculum.

- 92 1. Identifying species
- 93 2. Hoping that the noise averages out under large volume of data
- 94 3. Need for automated curation

95 Recommendations

- 96 1. Publish data (even the small one!)
- 97 2. Publish pipeline
- 98 3. Pay attention to standard when releasing data

99 **Acknowledgements** – This work was funded in part through a grant from the Canadian Institute of Ecology and
100 Evolution. TP was funded by a Starting grant from the Université de Montréal.

101 References

- 102 Barraquand, F. et al. 2014. Lack of quantitative training among early-career ecologists: a survey of the problem and
103 potential solutions. - PeerJ 2: e285.
- 104 Chamberlain, S. A. and Szöcs, E. 2013. taxize: taxonomic search and retrieval in R. - F1000Research in press.
- 105 Eklof, A. et al. 2011. Relevance of evolutionary history for food web structure. - Proc. R. Soc. B Biol. Sci. 279:
106 1588–1596.
- 107 Estes, J. A. et al. 2011. Trophic Downgrading of Planet Earth. - Science 333: 301–306.

108 Fensham, R. J. et al. 2011. Four desert waters: setting arid zone wetland conservation priorities through
109 understanding patterns of endemism. - Biol. Conserv. 144: 2459–2467.

110 Minckley, T. A. et al. 2013. The relevance of wetland conservation in arid regions: a re-examination of vanishing
111 communities in the American Southwest. - J. Arid Environ. 88: 213–221.

112 Poelen, J. H. et al. 2014. Global Biotic Interactions: An open infrastructure to share and analyze species-interaction
113 datasets. - Ecological Informatics in press.

114 Thuiller, W. et al. 2013. A road map for integrating eco-evolutionary processes into biodiversity models. - Ecol.
115 Lett. 16: 94–105.