

plantR: An R package and workflow for managing species records from biological collections

Renato A. F. de Lima^{1,2} | Andrea Sánchez-Tapia³ | Sara R. Mortara^{3,4} | Hans ter Steege^{1,5} | Martinez F. de Siqueira³

¹Tropical Botany, Naturalis Biodiversity Center, Leiden, The Netherlands

²Departamento de Ecología, Universidade de São Paulo, São Paulo, Brazil

³Instituto de Pesquisas Jardim Botânico do Rio de Janeiro, Rio de Janeiro, Brazil

⁴International Institute for Sustainability, Rio de Janeiro, Brazil

⁵Systems Ecology, Vrije Universiteit Amsterdam, Amsterdam, The Netherlands

Correspondence

Renato A. F. de Lima

Email: raflima@usp.br

Funding information

European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No 795114; Coordination for the Improvement of Higher Education Personnel - CAPES (process 88887.145924/2017-00);

Instituto Nacional da Mata Atlântica - INMA

Abstract

1. Species records from biological collections are becoming increasingly available online. This unprecedented availability of records has largely supported recent studies in taxonomy, biogeography, macroecology, and biodiversity conservation. Biological collections vary in their documentation and notation standards, which have changed through time. For different reasons, neither collections nor data repositories perform the editing, formatting, and standardization of the data, leaving these tasks to the final users of the species records (e.g. taxonomists, ecologists and conservationists). These tasks are challenging, particularly when working with millions of records from hundreds of biological collections.
2. To help collection curators and final users perform those tasks, we introduce `plantR`, an open-source package that provides a comprehensive toolbox to manage species records from biological collections. The package is accompanied by the proposal of a reproducible workflow to manage this type of data in taxonomy, ecology, and biodiversity conservation. It is implemented in R and designed to handle relatively large data sets as fast as possible. Initially designed to handle plant species records, many of the `plantR` features also apply to other groups of organisms, given that the data structure is similar.
3. The `plantR` workflow includes tools to (1) download records from different data repositories, (2) standardize typical fields associated with species records, (3) validate the locality, geographical coordinates, taxonomic nomenclature, and species identifications, including the retrieval of duplicates across collections, and (4) summarize and export records, including the construction of species checklists with vouchers.
4. Other R packages provide tools to tackle some of the workflow steps described above. But in addition to the new features and resources related to the data editing and validation, the greatest strength of `plantR` is to provide a comprehensive and user-friendly workflow in one single environment, performing all tasks from data retrieval to export. Thus, `plantR` can help researchers better assess data quality and avoid data leakage in a wide variety of studies using species records.

KEY WORDS

biodiversity, data cleaning, data download, duplicate records, gazetteer, GBIF, herbarium, taxonomic validation

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

1 | INTRODUCTION

2 Biological collections (e.g. museums and herbaria) are
3 essential for studying biodiversity (Graham et al., 2004).
4 Taxonomists use these collections to describe new
5 species, produce taxonomic revisions and species check-
6 lists, among other important uses (Funk, 2003; Bebber
7 et al., 2010; Besnard et al., 2018). In macroecology, bio-
8 geography, and conservation, biological collections are
9 often the main source of species records, which are used
10 to study spatial patterns of biodiversity, species ecolog-
11 ical niches, endemism levels, and conservation status
12 (Graham et al., 2004; Dauby et al., 2017; Ulloa et al.,
13 2017; Lima et al., 2020). Biological collections are in-
14 creasingly making their electronic databases available in
15 online databases, such as the Global Biodiversity Infor-
16 mation Facility (GBIF). This growing availability of infor-
17 mation has catalyzed many syntheses of our biodiversity
18 knowledge (e.g. Antonelli et al. 2018), highlighting the
19 importance of biological collections even more.

20 The increasing availability of biological collections
21 databases has also exposed the wide variation of the
22 documentation standards within and between collec-
23 tions (Willemse et al., 2008). Within collections, spec-
24 imens collected by different people or in different pe-
25 riods may vary in their notation standards. The inter-
26 national documentation standards themselves are con-
27 stantly evolving (www.tdwg.org/standards). Moreover,
28 older records tend to have less associated information
29 (e.g. missing geographical coordinates) and may contain
30 names of localities that no longer exist (i.e. changing to-
31 ponyms). Between collections, differences may emerge
32 from different choices of documentation standards, on
33 how to enter specimen information in the electronic
34 databases, and on which fields should be entered first
35 in the face of limited resources. The staff of biological
36 collections often have little time to update the informa-
37 tion that has been already entered in their databases or
38 to correct data entry errors (e.g. typographical errors).
39 These tasks become more challenging as the number of
40 records in the collection increases.

41 Despite the global efforts to standardize the docu-
42 mentation of biodiversity information (e.g. Darwin Core

43 standards), there is still much variation within fields as-
44 sociated with species records. This variation is likely to
45 remain for years to come because biological collections
46 are often underfunded, undervalued, and understaffed
47 (de Gasper et al., 2020). Online databases, such as GBIF,
48 gather, store, flag, and check some but not all the infor-
49 mation provided by the data providers. This means that,
50 although highly valuable, the available databases from
51 biological collections are not always ready for use (Peter-
52 son et al., 2018). So, the final users of species records
53 (e.g. taxonomists, ecologists, and conservationists) of-
54 ten have to decide between performing those pro-
55 cedures themselves or trusting the data available without
56 knowing exactly the level of data quality. This is prob-
57 lematic because variation in data quality can impact the
58 outcomes of studies in taxonomy, ecology, and conser-
59 vation (Graham et al., 2004; Zizka et al., 2019; Rodrigues
60 et al., 2020). Thus, we still need comprehensive and re-
61 producible tools to manage species records from biologi-
62 cal collections, particularly regarding notation standards,
63 species identifications, duplicate records, and fine-scale
64 validation of the geographical coordinates.

2 | OVERVIEW

65 We present `plantR`, a new R package for managing
66 species records from biological collections. As a general
67 approach, `plantR` does not edit the original infor-
68 mation; it stores the standardized information in new
69 columns to assist collection curators in comparing orig-
70 inal and edited information. Much of the new function-
71 alities depend on gazetteers, maps, lists of taxonomists,
72 and plant collections, which are provided with the pack-
73 age. As its name suggests, `plantR` was initially designed
74 to manage plant records from herbaria, with some func-
75 tionalities being currently exclusive to plants. However,
76 if the input data has the required fields and data for-
77 mat, many `plantR` features should work for any group
78 of organisms. `plantR` should interest taxonomists, bio-
79 geographers, ecologists, and conservationists, as well
80 as curators of biological collections. The package is
81 implemented in R (R Core Team, 2020) and details on

83 its implementation and functionalities can be found at
84 <https://github.com/LimaRAF/plantR>.

85 3 | THE PLANTR WORKFLOW

86 plantR is accompanied by the proposal of a workflow to
87 process the information associated with species records
88 (Fig. 1). Here, we present the steps of this workflow and
89 the main plantR features to apply it. They are presented
90 in the order that the workflow should be applied. This
91 order aims to maximize the edition and validation of the
92 available information, although many plantR functional-
93 ities work independently from the previous steps of the
94 workflow.

95 3.1 | Data entry

96 Users can download species records directly from R,
97 which is currently done from the Centro de Refer-
98 ência em Informação Ambiental (CRIA, www.cria.org.br) and GBIF (www.gbif.org), using functions
99 `rspecieslink()` and `rgbif2()`, respectively. The func-
100 tion `rgbif2()` performs a search based on scientific
101 names using the `rgbif` package, but with a stand-
102 arized output to enter the plantR workflow. The func-
103 tion `rspeciesLink()` is more flexible allowing the user
104 to search by scientific name or any other taxonomic
105 level, collection, and locality. Since these two sources of
106 species records return different fields, a function is pro-
107 vided to guarantee their correspondence with the DwC
108 standards (function `formatDwc()`). Users can also load
109 their own data, which can be converted to the Darwin
110 Core (DwC) standards (<https://dwc.tdwg.org>) using
111 the function `formatDwc()`. Alternatively, users can im-
112 port data from zipped DwC-Archive files from a local
113 directory or from a link for data download provided by
114 GBIF (function `readData()`).

116 3.2 | Data editing

117 Data standardization is particularly important when
118 combining records from multiple collections, because

they not always follow the same documentation stan-
dards. plantR provides tools to edit and standardize the
notation of the information associated with the records,
which are very important for validating locality informa-
tion, assessing the confidence level of species identifica-
tions and searching duplicate records across collections
(see 3.3 Data validation).

3.2.1 | People's names and collection 126 information

127 The first edits performed by plantR regards the name
128 of collector and identifiers, collector's number and
129 collection year (function `formatOcc()`). By default,
130 people's names are returned in the Biodiversity Infor-
131 mation Standards format ([www.tdwg.org/standards/](http://www.tdwg.org/standards/hispid3/)
132 `hispid3()`), which is: last name + comma + initials sep-
133 arated by points (e.g. Gentry, A.H.). Name formatting
134 takes into account generational suffixes (e.g. Junior),
135 prepositions (e.g. da, dos, von), compound last names
136 (e.g. Saint-Hilaire), some titles (e.g. Dr., Profa.) and mul-
137 tiple collector names. plantR also standardizes the col-
138 lection codes using a database of over 5000 plant col-
139 lection names and their respective Index Herbariorum
140 or Index Xylariorum codes (function `getCode()`).
141

3.2.2 | Locality and spatial information

142 One of the innovations of plantR is the standardiza-
143 tion of records' locality information (i.e the DwC fields
144 "country", "stateProvince", "municipality" and "locality";
145 function `formatLoc()`). For instance, names are trans-
146 formed to English (e.g. Brasil or Brésil become Brazil)
147 and their notation is standardized (e.g. BR or BRA be-
148 come Brazil). In the case of missing locality information,
149 plantR performs some text mining aiming to retrieve
150 them from other fields. To make sure that the original
151 or retrieved locality information does exist, the package
152 cross-checks the locality information of records with a
153 gazetteer (function `getLoc()`). This cross-checking is
154 based on a standard name-string that hierarchically com-
155 bines the locality information at the best resolution avail-
156 able, thus avoiding spurious matches of same locality
157

names in different countries or states/provinces (function `strLoc()`). The default `plantR` gazetteer currently contains entries at country level for all countries and at the lowest administrative level available at GDAM (<https://gadm.org>) for all Latin American countries and dependent territories (e.g. U.S. Virgin Islands). For Brazil, the gazetteer also contains information at the locality level (e.g. farms, forest fragments, parks). Most importantly, users can provide their regional or personal gazetteers.

The gazetteer includes some of the most common spelling variants and historical changes to locality names (currently biased for Brazil), which allows collection curators to trace back the most up-to-date locality names to improve their databases (function `getAdmin()`). Additionally, `plantR` assigns a geographical coordinate from the gazetteer to all valid localities (function `getCoord()`), which can be used as working coordinates in the case of missing or problematic original coordinates. Besides the automated assignment of missing coordinates, the package formats the original geographical coordinates to obtain non-zero, non-missing coordinates in decimal degrees (function `prepCoord()`).

3.2.3 | Taxonomic information

`plantR` offers tools to format scientific name notation (function `fixSpecies()`), such as the isolation and removal of taxonomic rank (e.g. var., subsp.) and name modifiers (e.g. cf., aff.), which is important for records containing more raw taxonomic information (e.g. morpho-species, incomplete identifications). The package also standardizes the name of botanical families, using a list of valid family names and synonyms from the APG IV for angiosperms (Chase et al., 2016) and PPG I for lycophytes and ferns (Schuettpelz et al. 2016; function `prepFamily()`). If the family name is not found in the list, a search for a valid family name is performed based on the genus. Finally, the package can replace synonyms, orthographic variants and typographical errors in species names (function `prepSpecies()`), which is performed using functions from the packages `Taxonstand` (Cayuela et al., 2021) and `flora` (Carvalho, 2020). These

packages perform exact and fuzzy name matching from The Plant List (www.theplantlist.org/) and the Brazilian Flora 2020 project (<http://floradobrasil.jbrj.gov.br/>), respectively.

3.3 | Data validation

3.3.1 | Locality and spatial information

`plantR` compares the precision of the original locality information with the one obtained by the cross-checking with a gazetteer (function `validateLoc()`). This comparison allows to flag possible typographical errors or unknown place names, which users can drop from the analyses or double-check themselves depending on their goals. Obtaining valid locality information is essential for the validation of geographical coordinates because they are validated by comparing the locality information of the record and the locality obtained by overlapping the coordinates with administrative maps (function `checkCoord()`). The package offers procedures for detecting the inversion and/or swap of coordinates (function `checkInverted()`), coordinates falling in the sea or bays, near the shoreline (`checkShore()`), and in neighbouring countries (`checkBorders()`). If after these procedures the locality information from the record and maps matches, the coordinate is flagged as validated, with an indication of the resolution of the validation (i.e. country, state, municipality or locality levels). As before, the validation of geographical coordinates is done using maps at the country level for the world and at the lowest administrative level available at GDAM for Latin America, but users can provide their own maps. Finally, `plantR` also provides tools to detect records from cultivated individuals (function `getCult()`) and spatial outliers (function `checkOut()`), i.e. coordinates too far away from the core distributions for a given taxon (Liu et al., 2018).

3.3.2 | Species identifications

One highlight of `plantR` is the classification of records according to the confidence in their species identifica-

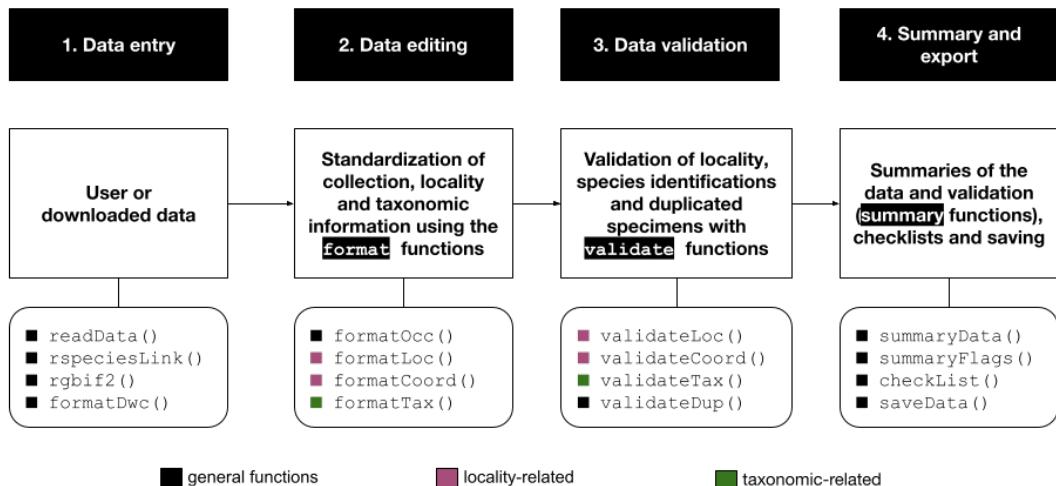


FIGURE 1 Chart illustrating the four main steps of the workflow proposed here to manage species records from biological collections for taxonomy, ecology, and biodiversity conservation. Black boxes represent each of the four steps, white boxes their description and rounded boxes their main *plantR* functions.

237 tions (function `validateTax()`). This validation is based
 238 on a global list of ca. 8500 plant taxonomists names
 239 compiled from different sources (Lima et al., 2020). By
 240 default, this classification assigns the highest confidence
 241 level to three different cases: (i) type specimens (e.g. iso-
 242 types, holotypes), (ii) records identified by a specialist
 243 of the family, and (iii) records collected by the special-
 244 ist of the family but with the identifier field empty (case
 245 iii is optional). The confidence level of records without
 246 identifier information (including NA's) is flagged as 'un-
 247 known', while records identified by non-family special-
 248 ists it are flagged as 'low'. Users can provide their own
 249 list of taxonomists, as long as this list has the same gen-
 250 eral format as the default list provided by *plantR*. More-
 251 over, `validateTax()` returns the most frequent names
 252 of identifiers that are not in the taxonomist list, allowing
 253 users to provide missing taxonomist names.

254 3.3.3 | Duplicate records

255 Another novelty of *plantR* regards duplicates, i.e. sam-
 256 ples of the same specimen incorporated in two or more
 257 collections (function `validateDup()`). Sharing biologi-

cal material across collections is a common and encour-
 258 aged practice, and they can represent 25% or more of
 259 the records available for regional biotas (e.g. Lima et al.,
 260 2020). The search for duplicates in *plantR* is executed
 261 by combining fields related to the taxonomy, collection
 262 and locality of the records (e.g., family + collector name
 263 + collector number + municipality). Because of the great
 264 variation in the notation and completeness of collec-
 265 tor's and localities names, the package allows the sim-
 266 ultaneous use of different combinations of these fields
 267 to search for duplicates (function `getDup()`). If two or
 268 more combinations are provided, the search of dupli-
 269 cates uses tools from network analysis to find both di-
 270 rect and indirect links between records. The retrieval
 271 of duplicates across collections performs well using rel-
 272 atively large data-sets (i.e. millions of records). How-
 273 ever, finding all existing duplicates requires that the
 274 databases of all collections are available and that all
 275 search fields are complete and filled in without typos
 276 using the same notation standards (or notations that
 277 *plantR* can standardize). This is rarely the case, so the
 278 list of duplicates returned should be considered incom-
 279 plete in many cases.

281 plantR provides not only tools to search for duplicates, but also to homogenize information within the
 282 groups of duplicates found, such as species, locality
 283 and/or spatial information (function `mergeDup()`). This
 284 homogenization allows retrieving the best information
 285 available within duplicates, which is particularly useful
 286 when collections vary in the number and completeness
 287 of the digitized fields. After this homogenization, users
 288 can choose to remove or not the duplicates from the
 289 data. See Lima et al. (2020) for more details on the
 290 search and merge of duplicates implemented here.

292 3.4 | Data summary and export

293 As a final step of the workflow, plantR can help users
 294 to summarize their data (e.g. number of occurrences,
 295 collections and species; function `summaryData()`) and
 296 the flags of the validation process (i.e. localities,
 297 coordinates, identifications and duplicates; function
 298 `summaryFlags()`). The package also provides species
 299 checklists with user-defined numbers of voucher speci-
 300 mens and the export of records by groups (e.g. families,
 301 countries, collections).

302 4 | IMPLEMENTATION

303 4.1 | Example of usage

304 The plantR workflow can be implemented using few
 305 command lines and wrapper functions (see Table 1 for
 306 details). Here, we provide a simple example using only
 307 one species. A detailed tutorial of the package is pro-
 308 vided at <https://github.com/LimaRAF/plantR>.

309
 310 # Installing plantR
 311 remotes::install_github("LimaRAF/plantR")
 312 library("plantR")
 313
 314 # Data download
 315 occs_splink <- rspeciesLink(species =
 316 "Euterpe edulis")
 317 occs_gbif <- rgbif2(species =
 318 "Euterpe edulis")

319
 320
 321
 322
 323
 324
 325
 326
 327
 328
 329
 330
 331
 332
 333
 334
 335
 336
 337
 338
 339
 340

occs <- formatDwc(splink_data =
 occs_splink,
 gbif_data =
 occs_gbif)

Data editing
 occs <- formatOcc(occs)
 occs <- formatLoc(occs)
 occs <- formatCoord(occs)
 occs <- formatTax(occs)

Data validation
 occs <- validateLoc(occs)
 occs <- validateCoord(occs)
 occs <- validateTax(occs)
 occs <- validateDup(occs)

Data summary
 summs <- summaryData(occs)
 flags <- summaryFlags(occs)
 checklist <- checkList(occs)

4.2 | Dependencies on other packages

341 Some of plantR's features depend on other R pack-
 342 ages (Table 1). Function `rgbif2()` uses package
 343 `rgbif` (Chamberlain et al., 2021) for downloading GBIF
 344 data. The management of strings, countries names,
 345 and spatial data use packages `stringr` (Wickham,
 346 2019), `countrycode` (Arel-Bundock et al., 2018), and
 347 `sf`, (Pebesma, 2018), respectively. As mentioned above,
 348 function `prepSpecies()` uses `Taxonstand` (Cayuela
 349 et al., 2021) and `flora` (Carvalho, 2020). The search
 350 of duplicates uses package `igraph` (Csardi and Nepusz,
 351 2006) to perform indirect string search. Finally, many
 352 functions use `data.table` (Dowle and Srinivasan, 2020),
 353 which provides fast table manipulation, reading and sav-
 354 ing.

TABLE 1 List of the main functions per type of information and per step of the proposed workflow. We also present the wrappers of the main functions for each step (if present) and the other R packages necessary to execute them.

Workflow step	Type of information	Main functions	Wrapper	Dependencies
1 - Data Entry	Species records	readData, rgbf2, rspeciesLink, formatDwc	-	rgbf, data.table
2 - Data Editing	Names, numbers, etc	prepName, colNumber, getYear, getCode	formatOcc	stringr
	Localities	fixLoc, strLoc, prepLoc, getLoc	formatLoc	countrycode, stringr
	Coordinates	prepCoord, getCoord	formatCoord	-
3 - Data Validation	Taxonomy	fixSpecies, prepSpecies, prepFamily	formatTax	flora, Taxonstand, data.table
	Localities	validateLoc	-	-
	Coordinates	checkCoord, checkBorders, checkShore, checkInverted, getCult, checkOut	validateCoord	sf, robustbase, data.table
4 - Summary and Export	Species identification	validateTax	-	-
	Duplicate records	prepDup, getDup, mergeDup, rmDup	validateDup	data.table, igraph
	Summaries	summaryData, summaryFlags, checklist	-	data.table, stringr
	Export	saveData	-	data.table

355 5 | DISCUSSION

356 5.1 | Comparison with other R packages

357 Other R packages already provide spelling and synonym
 358 checks of species names (Chamberlain and Szöcs 2013;
 359 Cayuela et al. 2021; Carvalho 2020; Kindt 2020), so
 360 there was no need to 'reinvent the wheel' and their
 361 functionalities were (or will be) integrated in plantR.
 362 CoordinateCleaner (Zizka et al., 2019) provides a great
 363 toolbox to work with geographical coordinates and we
 364 suggest this package for the advanced editing of geo-
 365 graphical coordinates. The differential of plantR lies
 366 in providing both locality and coordinate validation, the

367 automatic retrieval of coordinates for missing or prob-
 368 lematic coordinates, and the coordinate validation at
 369 the county level. However, because these validations
 370 depend on the package gazetteer, these innovations
 371 currently apply mainly to Latin America. plantR also
 372 provides an approach to find cultivated specimens (i.e.
 373 `getCult()`), which is based on the fields 'locality' or 'oc-
 374 currenceRemarks' and thus different from the approach
 375 used by `CoordinateCleaner`.

376 We found only one package that validates species
 377 identifications, `naturaList` (Rodrigues et al., 2020).
 378 This package also uses the field 'identifiedBy', but it re-
 379 turns more confidence levels of species identification

and requires a user-provided list of taxonomists. The differential of *plantR* relies on the provision of a large database of plant taxonomists, besides the possibility of the user providing an extra list of specialist names. In addition, *plantR* also relies on the field 'typeStatus' and it performs the validation at the family-level. We are not aware of other R packages that perform (i) the edition of people names, (ii) the validation of locality information and (iii) the search/merge of duplicates.

5.2 | Limitations and future developments

The variation in the notation of names, numbers and dates associated with species records across biological collections is huge; *plantR* handles most but not all of them. We envisage having a dictionary of common collectors' names, but today some double-checking is still necessary. As mentioned before, locality and county-level geographical validation are currently biased towards Latin America. Therefore, users must be aware that the package does not provide solutions to all problems related to species records information. Some improvements predicted to be implemented in the future include the download from other data repositories (e.g. JABOT, <http://jabot.jbrj.gov.br>), the expansion of the package gazetteer and county-level maps and the validation of species names against databases that have wider geographical and taxonomic coverage (e.g. Catalogue of Life). We also plan to include simple functions that prepare records to enter the workflow of other R packages (e.g. *modleR* or *ConR* - Sánchez-Tapia et al. 2020; Dauby et al. 2017), that facilitate the citation of collections (e.g. *occCite* - Owens et al. 2021) and that collect provenance (e.g. *rdt* - Lerner et al. 2018). Moreover, the gazetteer, list of taxonomists, maps, and collections are constantly being improved; we are open to receive and incorporate missing or regional information to make them more complete.

6 | CONCLUDING REMARKS

The number of collection databases made available online has greatly increased in the last decades and will probably continue to increase in the years to come (Graham et al., 2004; Sweeney et al., 2018). Therefore, having tools to assess and improve the quality of the information associated with species record is a pressing issue in biodiversity research. *plantR* provides these tools, some of them being presented for the first time. Although there are packages that provide similar tools, the greatest strength of *plantR* is to provide a comprehensive toolbox and a user-friendly workflow to process species records from beginning to end within a single environment. Thus, we expect that *plantR* can improve the reproducibility of taxonomic, ecological and conservation studies. But more importantly, we hope that *plantR* can assist collection curators to flag possible issues that need attention, thus saving their time while conducting the important task of maintaining biological collections.

ACKNOWLEDGEMENTS

This package was supported by the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No 795114. M.F.S., A.S.-T. and S.R.M. were supported by the Coordination for the Improvement of Higher Education Personnel - CAPES (process 88887.145924/2017-00), the PNPD/CAPES program and the PCI program of the 'Instituto Nacional da Mata Atlântica' (INMA), respectively. We thank Sidnei Souza from CRIA for his help with the web API. We also thank CNCFlora and the TreeCo database for providing localities used to construct the gazetteer, and Vinícius C. Souza (ESALQ/USP) who helped to curate the list of plant taxonomists.

AUTHORS' CONTRIBUTIONS

R.A.F.L. conceived the idea and R.A.F.L., A.S.-T., S.R.M. and M.F.S. designed methodology. R.A.F.L. constructed the list of taxonomists, collections, and families, while

455 R.A.F.L., A.S.-T., S.R.M. constructed the gazetteer and
 456 maps. R.A.F.L., A.S.-T., S.R.M. and H.t.S. wrote the codes
 457 and package documentations. R.A.F.L. led the writing of
 458 the manuscript, with contributions from A.S.-T. All au-
 459 thors contributed critically to the manuscript and gave
 460 final approval for publication.

461 DATA AVAILABILITY STATEMENT

462 The R package `plantR` is available at <https://github.com/LimaRAF/plantR>. The version of the package de-
 463 scribed in this paper (version 0.1.3) is archived at [link
 464 to be included before publication].

466 references

467 Antonelli, A., Zizka, A., Carvalho, F. A., Scharn, R., Bacon,
 468 C. D., Silvestro, D. and Condamine, F. L. (2018) Amazonia
 469 is the primary source of Neotropical biodiversity. *Proceedings of the National Academy of Sciences of the United
 470 States of America*, **115**, 6034–6039.

472 Arel-Bundock, V., Enevoldsen, N. and Yetman, C. (2018)
 473 countrycode: An R package to convert country names
 474 and country codes. *Journal of Open Source Software*, **3**,
 475 848. URL: <http://joss.theoj.org/papers/10.21105/joss.00848>.

477 Bebber, D. P., Carine, M. A., Wood, J. R. I., Wortley, A. H.,
 478 Harris, D. J., Prance, G. T., Davidse, G., Paige, J., Pen-
 479 nington, T. D., Robson, N. K. B. and Scotland, R. W.
 480 (2010) Herbaria are a major frontier for species discov-
 481 ery. *Proceedings of the National Academy of Sciences*,
 482 **107**, 22169–22171. URL: <http://www.pnas.org/cgi/doi/10.1073/pnas.1011841108>.

484 Besnard, G., Gaudeul, M., Lavergne, S., Muller, S., Rouhan,
 485 G., Sukhorukov, A. P., Vanderpoorten, A. and Jabbour,
 486 F. (2018) Herbarium-based science in the twenty-first
 487 century. *Botany Letters*, **165**, 323–327. URL: <https://doi.org/10.1080/23818107.2018.1482783>.

489 Carvalho, G. (2020) flora: Tools for Interacting with the
 490 Brazilian Flora 2020. *R package version 0.3.4*. URL:
 491 <https://cran.r-project.org/package=flora>.

492 Cayuela, L., Stein, A. and Oksanen, J. (2021) Taxonstand:
 493 Taxonomic Standardization of Plant Species Names. *R
 494 package version 2.3*. URL: <https://cran.r-project.org/package=Taxonstand>.

496 Chamberlain, S., Barve, V., Mcglinn, D., Oldoni, D., Desmet,
 497 P., Geffert, L. and Ram, K. (2021) `rgbif`: Interface to
 498 the Global Biodiversity Information Facility API. *R pack-
 499 age version 3.5.2*. URL: <https://cran.r-project.org/package=rgbif>.

500 Chamberlain, S. A. and Szöcs, E. (2013) `taxize`: taxonomic
 501 search and retrieval in R. *F1000Research*, **2**, 191. URL:
 502 <https://f1000research.com/articles/2-191/v1>.

503 Chase, M. W., Christenhusz, M. J., Fay, M. F., Byng, J. W.,
 504 Judd, W. S., Soltis, D. E., Mabberley, D. J., Sennikov, A. N.,
 505 Soltis, P. S., Stevens, P. F., Briggs, B., Brockington, S.,
 506 Chautems, A., Clark, J. C., Conran, J., Haston, E., Möller,
 507 M., Moore, M., Olmstead, R., Perret, M. et al. (2016) An
 508 update of the Angiosperm Phylogeny Group classifica-
 509 tion for the orders and families of flowering plants: APG
 510 IV. *Botanical Journal of the Linnean Society*, **181**, 1–20.
 511 URL: <https://academic.oup.com/botlinnean/article-lookup/doi/10.1111/bj.12385>.

512 Csardi, G. and Nepusz, T. (2006) The `igraph` software pack-
 513 age for complex network research. *InterJournal Complex
 514 Systems*, 1695. URL: <https://igraph.org>.

515 Dauby, G., Stévert, T., Droissart, V., Cosiaux, A., Deblauwe,
 516 V., Simo-Droissart, M., Sosef, M. S., Lowry, P. P., Schatz,
 517 G. E., Gereau, R. E. and Couvreur, T. L. (2017) ConR: An
 518 R package to assist large-scale multispecies preliminary
 519 conservation assessments using distribution data. *Ecol-
 520 ogy and Evolution*, **7**, 11292–11303.

521 Dowle, M. and Srinivasan, A. (2020) `data.table`: Extension
 522 of 'data.frame'. *R Package Version 1.13.6*. URL: <https://cran.r-project.org/package=data.table>.

523 Funk, V. (2003) The Importance of Herbaria. *Plant Science
 524 Bulletin*, **49**, 94–95.

525 de Gasper, A. L., Stehmann, J. R., Roque, N., Bigio, N. C.,
 526 Sartori, Á. L. B. and Gritt, G. S. (2020) Brazilian herbaria:
 527 An overview. *Acta Botanica Brasilica*, **34**, 352–359.

528 Graham, C., Ferrier, S., Huettman, F., Moritz, C. and Peter-
 529 son, A. (2004) New developments in museum-based
 530 informatics and applications in biodiversity analy-
 531 sis. *Trends in Ecology Evolution*, **19**, 497–503. URL:
 532 <https://linkinghub.elsevier.com/retrieve/pii/S0169534704002034>.

533 Kindt, R. (2020) `WorldFlora`: An R package for exact and
 534 fuzzy matching of plant names against the World Flora
 535 Online taxonomic backbone data. *Applications in Plant
 536 Sciences*, **8**. URL: <https://onlinelibrary.wiley.com/doi/10.1002/aps3.11388>.

542 Lerner, B., Boose, E. and Perez, L. (2018) Using Introspection 587
 543 to Collect Provenance in R. *Informatics*, **5**. URL: <http://www.mdpi.com/2227-9709/5/1/12>. 588
 544

545 Lima, R. A. F., Souza, V. C., de Siqueira, M. F. and ter Steege, 589
 546 H. (2020) Defining endemism levels for biodiversity 588
 547 conservation: Tree species in the Atlantic Forest hotspot. *Bi- 589
 548 ological Conservation*, **252**, 108825. URL: <https://doi.org/10.1016/j.biocon.2020.108825>. 590

550 Liu, C., White, M. and Newell, G. (2018) Detecting outliers 590
 551 in species distribution data. *Journal of Biogeography*, **45**, 591
 552 164–176. 592

553 Owens, H. L., Merow, C., Maitner, B., Kass, J. M., Barve, V. 593
 554 and Guralnick, R. P. (2021) occCite: Querying and 593
 555 Managing Large Biodiversity Occurrence Datasets. *R package 594
 556 version 0.4.6*. URL: <https://cran.r-project.org/package=occCite>. 595

558 Pebesma, E. (2018) Simple Features for R: Standardized 599
 559 Support for Spatial Vector Data. *The R Journal*, **10**, 600
 560 439. URL: <https://journal.r-project.org/archive/2018/RJ-2018-009/index.html>. 601

562 Peterson, A. T., Asase, A., Canhos, D., de Souza, S. and Wieczorek, J. (2018) Data Leakage and Loss in Biodiversity 602
 563 Informatics. *Biodiversity Data Journal*, **6**, e26826. URL: 603
 564 <https://bdj.pensoft.net/articles.php?id=26826>. 604

566 R Core Team (2020) R: A language and environment for 604
 567 statistical computing. *R Foundation for Statistical Computing*, 605
 568 Vienna, Austria. URL: <https://www.r-project.org>. 606

569 Rodrigues, A. V., Nakamura, G. and Duarte, L. (2020) natu- 606
 570 raList : a package to classify occurrence records in 607
 571 levels of confidence in species identification. *bioRxiv*, 1–17. 608
 572 URL: <https://doi.org/10.1101/2020.05.26.115220>. 609

573 Sánchez-Tapia, A., Mortara, S. R., Bezerra Rocha, D. S., 610
 574 Mendes Barros, F. S., Gall, G. and de Siqueira, M. F. 611
 575 (2020) modleR: a modular workflow to perform ecological 611
 576 niche modeling in R. *bioRxiv*, 1–25. 612

577 Schuettpelz, E., Schneider, H., Smith, A. R., Hovenkamp, P., 612
 578 Prado, J., Rouhan, G., Salino, A., Sundue, M., Almeida, 612
 579 T. E., Parris, B., Sessa, E. B., Field, A. R., de Gasper, A. L., 612
 580 Rothfels, C. J., Windham, M. D., Lehnert, M., Dauphin, 612
 581 B., Ebihara, A., Lehtonen, S., Schwartsburd, P. B. et al. 612
 582 (2016) A community-derived classification for extant lycophytes 612
 583 and ferns. *Journal of Systematics and Evolution*, 612
 584 **54**, 563–603.

585 Sweeney, P. W., Starly, B., Morris, P. J., Xu, Y., Jones, A., 612
 586 Radhakrishnan, S., Grassa, C. J. and Davis, C. C. (2018) 612

Large-scale digitization of herbarium specimens: Development and usage of an automated, high-throughput conveyor system. *Taxon*, **67**, 165–178.

Ulloa, C. U., Acevedo-Rodríguez, P., Beck, S., Belgrano, 590
 M. J., Bernal, R., Berry, P. E., Brako, L., Celis, M., 591
 Davidse, G., Forzza, R. C., Gradstein, S. R., Hokche, 592
 O., León, B., León-Yáñez, S., Magill, R. E., Neill, D. A., 593
 Nee, M., Raven, P. H., Stimmel, H., Strong, M. T. et al. (2017) An integrated assessment of the vascular 594
 plant species of the Americas. *Science*, **358**, 1614– 595
 1617. URL: <http://www.science.org/lookup/doi/10.1126/science.aao0398>.

Wickham, H. (2019) stringr: Simple, Consistent Wrappers 599
 for Common String Operations. *R package version 1.4.0*. 600
 URL: <https://cran.r-project.org/package=stringr>. 601

Willemse, L. P., Van Welzen, P. C. and Mols, J. B. (2008) 602
 Standardisation in data-entry across databases: Avoiding 603
 Babylonian confusion. *Taxon*, **57**, 343–345. 604

Zizka, A., Silvestro, D., Andermann, T., Azevedo, J., Duarte 605
 Ritter, C., Edler, D., Farooq, H., Herdean, A., Ariza, 606
 M., Scharn, R., Svantesson, S., Wengström, N., Zizka, 607
 V. and Antonelli, A. (2019) CoordinateCleaner: Stan- 608
 dardized cleaning of occurrence records from biological 609
 collection databases. *Methods in Ecology and Evolution*, 610
10, 744–751. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1111/2041-210X.13152>.