

# Package ‘naturalist’

October 13, 2022

**Type** Package

**Title** Classify Occurrences by Confidence Levels in the Species ID

**Version** 0.5.0

**Description** Classify occurrence records based on confidence levels of species identification. In addition, implement tools to filter occurrences inside grid cells and to manually check for possible errors with an interactive shiny application.

**License** MIT + file LICENSE

**Encoding** UTF-8

**LazyData** true

**RoxygenNote** 7.1.2

**Imports** shiny, shinyWidgets, dplyr, stringr, sp, raster, shinydashboard, leaflet, leaflet.extras, tidytext, magrittr, vegan, fasterize, sf, htmltools, methods, rlang

**Suggests** knitr, rmarkdown, testthat, tm, rgeos, rgdal, rnaturalearth, lwgeom, shinyLP, maptools

**VignetteBuilder** knitr

**Depends** R (>= 2.10)

**URL** <https://github.com/avrodrigues/naturalist>

**BugReports** <https://github.com/avrodrigues/naturalist/issues>

**NeedsCompilation** no

**Author** Arthur Vinicius Rodrigues [aut, cre] (<<https://orcid.org/0000-0003-2656-558X>>), Gabriel Nakamura [aut] (<<https://orcid.org/0000-0002-5144-5312>>), Leandro Duarte [aut] (<<https://orcid.org/0000-0003-1771-0407>>)

**Maintainer** Arthur Vinicius Rodrigues <rodrigues.arthur.v@gmail.com>

**Repository** CRAN

**Date/Publication** 2022-04-20 13:30:02 UTC

## R topics documented:

A.setosa . . . . .	2
BR . . . . .	3
classify_occ . . . . .	3
clean_eval . . . . .	6
create_spec_df . . . . .	9
cyathea.br . . . . .	10
define_env_space . . . . .	10
get_det_names . . . . .	11
grid_filter . . . . .	12
map_module . . . . .	14
r.temp.prec . . . . .	17
specialists . . . . .	17
spec_names_ex . . . . .	18

<b>Index</b>	<b>19</b>
--------------	-----------

---

A.setosa	<i>Occurrence records of Alsophila setosa downloaded from Global Biodiversity Information Facility (GBIF).</i>
----------	--

---

### Description

A GBIF raw dataset containing 508 occurrence records for the tree fern *Alsophila setosa*.

### Usage

```
A.setosa
```

### Format

A data frame with 508 rows and 45 variables

### Source

GBIF.org (08 July 2019) GBIF Occurrence Download doi: [10.15468/dl.6jesg0](https://doi.org/10.15468/dl.6jesg0)

---

BR	<i>Brazil boundary</i>
----	------------------------

---

**Description**

A spatial polygon with the Brazil boundaries

**Usage**

BR

**Format**

A 'SpatialPolygonsDataFrame' with 1 feature

---

classify_occ	<i>Classify occurrence records in levels of confidence in species identification</i>
--------------	--

---

**Description**

Classifies occurrence records in levels of confidence in species identification

**Usage**

```
classify_occ(
  occ,
  spec = NULL,
  na.rm.coords = TRUE,
  crit.levels = c("det_by_spec", "not_spec_name", "image", "sci_collection",
    "field_obs", "no_criteria_met"),
  ignore.det.names = NULL,
  spec.ambiguity = "not.spec",
  institution.code = "institutionCode",
  collection.code = "collectionCode",
  catalog.number = "catalogNumber",
  year = "year",
  date.identified = "dateIdentified",
  species = "species",
  identified.by = "identifiedBy",
  decimal.latitude = "decimalLatitude",
  decimal.longitude = "decimalLongitude",
  basis.of.record = "basisOfRecord",
  media.type = "mediaType",
  occurrence.id = "occurrenceID",
```

```

institution.source,
year.event,
scientific.name,
determined.by,
latitude,
longitude,
basis.of.rec,
occ.id
)

```

### Arguments

<code>occ</code>	data frame with occurrence records information.
<code>spec</code>	data frame with specialists' names. See details.
<code>na.rm.coords</code>	logical. If TRUE, remove occurrences with NA in <code>decimal.latitude</code> or <code>decimal.longitude</code>
<code>crit.levels</code>	character. Vector with levels of confidence in decreasing order. The criteria allowed are <code>det_by_spec</code> , <code>not_spec_name</code> , <code>image</code> , <code>sci_collection</code> , <code>field_obs</code> , <code>no_criteria_met</code> . See details.
<code>ignore.det.names</code>	character vector indicating strings in <code>identified.by</code> that should be ignored as a taxonomist. See details.
<code>spec.ambiguity</code>	character. Indicates how to deal with ambiguity in specialists names. <code>not.spec</code> solve ambiguity by classifying the identification as done by a non-specialist; <code>is.spec</code> assumes the identification was done by a specialist; <code>manual.check</code> enables the user to manually check all ambiguous names. Default is <code>not.spec</code> .
<code>institution.code</code>	column name of <code>occ</code> with the name (or acronym) in use by the institution having custody of the object(s) or information referred to in the record.
<code>collection.code</code>	column name of <code>occ</code> with The name, acronym, code, or initials identifying the collection or data set from which the record was derived.
<code>catalog.number</code>	column name of <code>occ</code> with an identifier (preferably unique) for the record within the data set or collection.
<code>year</code>	Column name of <code>occ</code> the four-digit year in which the Event occurred, according to the Common Era Calendar.
<code>date.identified</code>	Column name of <code>occ</code> with the date on which the subject was determined as representing the Taxon.
<code>species</code>	column name of <code>occ</code> with the species names.
<code>identified.by</code>	column name of <code>occ</code> with the name of who determined the species.
<code>decimal.latitude</code>	column name of <code>occ</code> latitude in decimal degrees.
<code>decimal.longitude</code>	column name of <code>occ</code> longitude in decimal degrees.
<code>basis.of.record</code>	column name with the specific nature of the data record. See details.

media.type	column name of occ with the media type of recording. See details.
occurrence.id	column name of occ with link or code for the occurrence record. See in <a href="#">Darwin Core Format</a>
institution.source	deprecated, use institution.code instead.
year.event	deprecated, use year instead.
scientific.name	deprecated, use species instead.
determined.by	deprecated, use identified.by instead
latitude	deprecated, use decimal.latitude instead
longitude	deprecated, use decimal.longitude instead
basis.of.rec	deprecated, use basis.of.record instead.
occ.id	deprecated, use occurrence.id instead

### Details

spec data frame must have columns separating LastName, Name and Abbrev. See [create\\_spec\\_df](#) function for a easy way to produce this data frame.

When ignore.det.name = NULL (default), the function ignores strings with "RRC ID Flag", "NA", "", "-\_" and "\_". When a character vector is provided, the function adds the default strings to the provided character vector and ignore all these strings as being a name of a taxonomist.

The function classifies the occurrence records in six levels of confidence in species identification. The six levels are:

- det\_by\_spec - when the identification was made by a specialists which is present in the list of specialists provided in the spec argument;
- not\_spec\_name - when the identification was made by a name who is not a specialist name provide in spec;
- image - the occurrence have not name of a identifier, but present an image associated;
- sci\_collection - the occurrence have not name of a identifier, but preserved in a scientific collection;
- field\_obs - the occurrence have not name of a identifier, but it was identified in field observation;
- no\_criteria\_met - no other criteria was met.

The (decreasing) order of the levels in the character vector determines the classification level order.

basis.of.record is a character vector with one of the following types of record: PRESERVED\_SPECIMEN, PreservedSpecimen, HUMAN\_OBSERVATION or HumanObservation, as in GBIF data 'basisOfRecord'.

media.type uses the same pattern as GBIF mediaType column, indicating the existence of an associated image with stillImage.

### Value

The occ data frame plus the classification of each record in a new column, named naturalist\_levels.

**Author(s)**

Arthur V. Rodrigues

**See Also**

[specialists](#)

**Examples**

```
data("A.setosa")
data("specialists")
occ.class <- classify_occ(A.setosa, specialists)
```

---

clean\_eval

*Evaluate the cleaning of occurrences records*

---

**Description**

This function compare the area occupied by a species before and after pass through the cleaning procedure according to the chosen level of filter. The comparison can be made by measuring area in the geographical and in the environmental space

**Usage**

```
clean_eval(
  occ.cl,
  geo.space,
  env.space = NULL,
  level.filter = c("1_det_by_spec"),
  r,
  species = "species",
  decimal.longitude = "decimalLongitude",
  decimal.latitude = "decimalLatitude",
  scientific.name,
  longitude,
  latitude
)
```

**Arguments**

occ.cl	data frame with occurrence records information already classified by <a href="#">classify_occ</a> function.
geo.space	a SpatialPolygons* or sf object defining the geographical space

<code>env.space</code>	a <code>SpatialPolygons*</code> or <code>sf</code> object defining the environmental space. Use the <a href="#">define_env_space</a> for create this object. By default <code>env.space = NULL</code> , hence do not evaluate the cleaning in the environmental space.
<code>level.filter</code>	a character vector including the levels in 'naturaList_levels' column which filter the occurrence data set.
<code>r</code>	a raster with 2 layers representing the environmental variables. If <code>env.space = NULL</code> , it could be a single layer raster, from which the cell size and extent are extracted to produce the composition matrix.
<code>species</code>	column name of <code>occ.cl</code> with the species names.
<code>decimal.longitude</code>	column name of <code>occ.cl</code> longitude in decimal degrees.
<code>decimal.latitude</code>	column name of <code>occ.cl</code> latitude in decimal degrees.
<code>scientific.name</code>	deprecated, use <code>species</code> instead.
<code>longitude</code>	deprecated, use <code>decimal.longitude</code> instead
<code>latitude</code>	deprecated, use <code>decimal.latitude</code> instead

### Value

a list in which:

`area` data frame remaining area after cleaning proportional to the area before cleaning. The values vary from 0 to 1. Column named `r.geo.area` is the remaining area for all species in the geographic space and the `r.env.area` in the environmental space.

`comp` data frame with composition of species in sites (cells from raster layers) before cleaning (`comp$comp$BC`) and after cleaning (`comp$comp$AC`). The number of rows is equal the number of cells in `r`, and number of columns is equal to the number of species in the `occ.cl`.

`rich` data frame with a single column with the richness of each site

`site.coords` data frame with site's coordinates. It facilitates to built raster layers from results using [rasterFromXYZ](#)

### See Also

[define\\_env\\_space](#)

### Examples

```
## Not run:

library(sp)
library(raster)

data("specialists") # list of specialists
data("cyathea.br") # occurrence dataset
```

```

# classify
occ.cl <- classify_occ(cyathea.br, specialLists)

# delimit the geographic space
# land area
data("BR")

# Transform occurrence data in SpatialPointsDataFrame
spdf.occ.cl <- sp::SpatialPoints(occ.cl[, c("decimalLongitude", "decimalLatitude")])

# load climate data
data("r.temp.prec") # mean temperature and annual precipitation
df.temp.prec <- raster::as.data.frame(r.temp.prec)

### Define the environmental space for analysis
# this function will create a boundary of available environmental space,
# analogous to the continent boundary in the geographical space
env.space <- define_env_space(df.temp.prec, buffer.size = 0.05)

# filter by year to be consistent with the environmental data
occ.class.1970 <- occ.cl %>%
  dplyr::filter(year >= 1970)

### run the evaluation
cl.eval <- clean_eval(occ.class.1970,
                     env.space = env.space,
                     geo.space = BR,
                     r = r.temp.prec)

#area results
head(cl.eval$area)

### richness maps
## it makes sense if there are more than one species
rich.before.clean <- raster::rasterFromXYZ(cbind(cl.eval$site.coords,
                                                  cl.eval$rich$rich.BC))
rich.after.clean <- raster::rasterFromXYZ(cbind(cl.eval$site.coords,
                                                  cl.eval$rich$rich.AC))

raster::plot(rich.before.clean)
raster::plot(rich.after.clean)

### species area map
comp.bc <- as.data.frame(cl.eval$comp$comp.BC)
comp.ac <- as.data.frame(cl.eval$comp$comp.AC)

c.villosa.bc <- raster::rasterFromXYZ(cbind(cl.eval$site.coords,
                                             comp.bc`Cyathea villosa`))
c.villosa.ac <- raster::rasterFromXYZ(cbind(cl.eval$site.coords,

```



```
comp.ac$`Cyathea villosa`))

raster::plot(c.villosa.bc)
raster::plot(c.villosa.ac)

## End(Not run)
```

---

create_spec_df	<i>Create specialist data frame from character vector</i>
----------------	---

---

## Description

Creates a specialist data frame ready for use in [classify\\_occ](#) from a character vector containing the specialists names

## Usage

```
create_spec_df(spec.char)
```

## Arguments

spec.char      a character vector with specialist names

## Value

a data frame. Columns split the names, surname and abbreviation for the names. If the full name contain any special character, such as accent marks, two lines for that name will be provided, with and without the special characters. See examples.

## Examples

```
# Example using Latin accent marks
data(spec_names_ex)

spec_names_ex
create_spec_df(spec_names_ex)
```

---

cyathea.br	<i>Occurrence records of Cyathea species in Brazil downloaded from Global Biodiversity Information Facility (GBIF).</i>
------------	---

---

### Description

A filtered GBIF dataset containing 3851 occurrence records for the fern species from the genus *Cyathea* in Brazil. We filtered the data after download from GBIF to ensure all occurrences records are from Brazil.

### Usage

```
cyathea.br
```

### Format

A data frame with 3851 rows and 50 variables

### Source

GBIF.org (07 March 2021) GBIF Occurrence Download doi: [10.15468/dl.qrhynv](https://doi.org/10.15468/dl.qrhynv)

---

define_env_space	<i>Define environmental space for species occurrence</i>
------------------	--

---

### Description

Based on two continuous environmental variables, it defines a bi-dimensional environmental space.

### Usage

```
define_env_space(env, buffer.size, plot = TRUE)
```

### Arguments

env	matrix or data frame with two columns containing two environmental variables. The variables must be numeric, even for data frames.
buffer.size	numeric value indicating a buffer size around each point which will delimit the environmental geographical border for the occurrence point. See details.
plot	logical. whether to plot the polygon. Default is TRUE.

**Details**

The environmental variables are standardized by range, which turns the range of each environmental variable from 0 to 1. Then, it is delimited a buffer of size equal to `buffer.size` around each point in this space and a polygon is draw to link these buffers. The function returns the polygon needed to link all points, and the area of the polygon indicates the environmental space based in the variables used.

**Value**

An object of `sfc_POLYGON` class

**Examples**

```
## Not run:
library("raster")

# load climate data
data("r.temp.prec")
env.data <- raster::as.data.frame(r.temp.prec)

define_env_space(env.data, 0.05)

## End(Not run)
```

---

`get_det_names`*Get the names in the 'identified.by' column*

---

**Description**

This function facilitates the search for non-taxonomist strings in the 'identified.by' column of occurrence records data set

**Usage**

```
get_det_names(  
  occ,  
  identified.by = "identifiedBy",  
  freq = FALSE,  
  decreasing = TRUE,  
  determined.by  
)
```

**Arguments**

occ	data frame with occurrence records information.
identified.by	column name of occ with the name of who determined the species.
freq	logical. If TRUE output contain the number of times each string is repeated in the identified.by column. Default = FALSE
decreasing	logical. sort strings in decreasing order of frequency. Default = TRUE.
determined.by	deprecated, use identified.by instead.

**Value**

character vector containing the strings in identified.by column of occ. If freq = TRUE it return a data frame with two columns: 'strings' and 'frequency'.

**Examples**

```
data("A.setosa")
get_det_names(A.setosa, freq = TRUE)
```

---

grid_filter	<i>Filter the occurrence with most confidence in species identification inside grid cells</i>
-------------	---

---

**Description**

In each grid cell it selects the occurrence with the highest confidence level in species identification made by `classify_occ` function.

**Usage**

```
grid_filter(
  occ.cl,
  grid.resolution = c(0.5, 0.5),
  r = NULL,
  institution.code = "institutionCode",
  collection.code = "collectionCode",
  catalog.number = "catalogNumber",
  year = "year",
  date.identified = "dateIdentified",
  species = "species",
  identified.by = "identifiedBy",
  decimal.latitude = "decimalLatitude",
  decimal.longitude = "decimalLongitude",
  basis.of.record = "basisOfRecord",
  media.type = "mediaType",
  occurrence.id = "occurrenceID",
```

```

institution.source,
year.event,
scientific.name,
determined.by,
latitude,
longitude,
basis.of.rec,
occ.id
)

```

## Arguments

occ.cl	data frame with occurrence records information already classified by <code>classify_occ</code> function.
grid.resolution	numeric vector with width and height of grid cell in decimal degrees.
r	raster from which the grid cell resolution is derived.
institution.code	column name of <code>occ.cl</code> with the name (or acronym) in use by the institution having custody of the object(s) or information referred to in the record.
collection.code	column name of <code>occ.cl</code> with The name, acronym, code, or initials identifying the collection or data set from which the record was derived.
catalog.number	column name of <code>occ.cl</code> with an identifier (preferably unique) for the record within the data set or collection.
year	Column name of <code>occ.cl</code> the four-digit year in which the Event occurred, according to the Common Era Calendar.
date.identified	Column name of <code>occ.cl</code> with the date on which the subject was determined as representing the Taxon.
species	column name of <code>occ</code> with the species names.
identified.by	column name of <code>occ.cl</code> with the name of who determined the species.
decimal.latitude	column name of <code>occ.cl</code> latitude in decimal degrees.
decimal.longitude	column name of <code>occ.cl</code> longitude in decimal degrees.
basis.of.record	column name with the specific nature of the data record. See details.
media.type	column name of <code>occ.cl</code> with the media type of recording. See details.
occurrence.id	column name of <code>occ</code> with link or code for the occurrence record. See in <a href="#">Darwin Core Format</a>
institution.source	deprecated, use <code>institution.code</code> instead.
year.event	deprecated, use <code>year</code> instead.

scientific.name	deprecated, use species instead.
determined.by	deprecated, use identified.by instead
latitude	deprecated, use decimal.latitude instead
longitude	deprecated, use decimal.longitude instead
basis.of.rec	deprecated, use basis.of.record instead.
occ.id	deprecated, use occurrence.id instead

**Value**

Data frame with the same columns of `occ.cl`.

**Author(s)**

Arthur V. Rodrigues

**See Also**

[classify\\_occ](#)

**Examples**

```
## Not run:  
  
data("A.setosa")  
data("specialists")  
  
occ.class <- classify_occ(A.setosa, specialists)  
occ.grid <- grid_filter(occ.class)  
  
## End(Not run)
```

---

map\_module

*Check the occurrence records in a interactive map module*

---

**Description**

Allows to delete occurrence records and to select occurrence points by classification levels or by drawing spatial polygons.

**Usage**

```

map_module(
  occ.cl,
  action = "clean",
  institution.code = "institutionCode",
  collection.code = "collectionCode",
  catalog.number = "catalogNumber",
  year = "year",
  date.identified = "dateIdentified",
  species = "species",
  identified.by = "identifiedBy",
  decimal.latitude = "decimalLatitude",
  decimal.longitude = "decimalLongitude",
  basis.of.record = "basisOfRecord",
  media.type = "mediaType",
  occurrence.id = "occurrenceID",
  institution.source,
  year.event,
  scientific.name,
  determined.by,
  latitude,
  longitude,
  basis.of.rec,
  occ.id
)

```

**Arguments**

occ.cl	Data frame with occurrence records information already classified by <a href="#">classify_occ</a> function.
action	a string with "clean" or "flag" which defines the action of 'map_module' function with the occurrence dataset. Default is "clean". If the string is "clean" the dataset returned only the occurrences records selected by the user. If the string is "flag", a column named 'map_module_flag' is added in the output dataset, with tags 'selected' and 'deleted', following the choices of the user in the application.
institution.code	column name of occ with the name (or acronym) in use by the institution having custody of the object(s) or information referred to in the record.
collection.code	column name of occ with The name, acronym, code, or initials identifying the collection or data set from which the record was derived.
catalog.number	column name of occ with an identifier (preferably unique) for the record within the data set or collection.
year	Column name of occ the four-digit year in which the Event occurred, according to the Common Era Calendar.

date.identified	Column name of occ with the date on which the subject was determined as representing the Taxon.
species	column name of occ with the species names.
identified.by	column name of occ with the name of who determined the species.
decimal.latitude	column name of occ latitude in decimal degrees.
decimal.longitude	column name of occ longitude in decimal degrees.
basis.of.record	column name with the specific nature of the data record. See details.
media.type	column name of occ with the media type of recording. See details.
occurrence.id	column name of occ with link or code for the occurrence record. See in <a href="#">Darwin Core Format</a>
institution.source	deprecated, use institution.code instead.
year.event	deprecated, use year instead.
scientific.name	deprecated, use species instead.
determined.by	deprecated, use identified.by instead
latitude	deprecated, use decimal.latitude instead
longitude	deprecated, use decimal.longitude instead
basis.of.rec	deprecated, use basis.of.record instead.
occ.id	deprecated, use occurrence.id instead

**Value**

Data frame with the same columns of occ.cl.

**Author(s)**

Arthur V. Rodrigues

**See Also**

[classify\\_occ](#)

**Examples**

```
## Not run:
data("A.setosa")
data("specialists")

occ.class <- classify_occ(A.setosa, specialists)
occ.selected <- map_module(occ.class)
occ.selected
```



```
## End(Not run)
```

---

```
r.temp.prec          Raster of temperature and precipitation
```

---

### Description

Raster of Annual Mean Temperature (bio1) and Total Annual Precipitation (bio2). Layers were downloaded from worldclim database and cropped to the extent of cyathea\_br with a buffer of 100 km.

### Usage

```
r.temp.prec
```

### Format

A raster with two layers

---

```
specialists          Specialists of ferns and lycophytes of Brazil
```

---

### Description

A dataset containing the specialists of ferns and lycophytes of Brazil formatted to be used by naturalist package. This data serves as a format example for spec argument in [classify\\_occ](#).

### Usage

```
specialists
```

### Format

A data frame with 27 rows and 8 columns:

**LastName** Last name of the specialist.

**Name1** Columns with the names of specialist. Could be repeated as long as needed. In this data Name\* was repeated three times.

**Name2** Columns with the names of specialist.

**Name3** Columns with the names of specialist.

**Name4** Columns with the names of specialist.

**Abbrev1** Columns with the abbreviation (one character) of the names of specialists. Could be repeated as long as needed. In this data Abbrev\* was repeated three times.

**Abbrev2** Columns with the abbreviation (one character) of the names of specialists.

**Abbrev3** Columns with the abbreviation (one character) of the names of specialists.

**Source**

The specialists names was derived from the authors of paper: doi: [10.1590/21757860201566410](https://doi.org/10.1590/21757860201566410)

---

spec\_names\_ex

*Example of specialist names with accent marks*

---

**Description**

Example of specialist names with accent marks

**Usage**

spec\_names\_ex

**Format**

character

# Index

## \* datasets

- A.setosa, [2](#)
- BR, [3](#)
- cyathea.br, [10](#)
- r.temp.prec, [17](#)
- spec\_names\_ex, [18](#)
- specialists, [17](#)

A.setosa, [2](#)

BR, [3](#)

classify\_occ, [3](#), [6](#), [9](#), [12–17](#)

clean\_eval, [6](#)

create\_spec\_df, [5](#), [9](#)

cyathea.br, [10](#)

define\_env\_space, [7](#), [10](#)

get\_det\_names, [11](#)

grid\_filter, [12](#)

map\_module, [14](#)

r.temp.prec, [17](#)

rasterFromXYZ, [7](#)

spec\_names\_ex, [18](#)

specialists, [6](#), [17](#)