Species occurrence data (spocc), version 0.2.0

Introduction

The rOpenSci projects aims to provide programmatic access to scientific data repositories on the web. A vast majority of the packages in our current suite retrieve some form of biodiversity or taxonomic data. Since several of these datasets have been georeferenced, it provides numerous opportunities for visualizing species distributions, building species distribution maps, and for using it analyses such as species distribution models. In an effort to streamline access to these data, we have developed a package called Spocc, which provides a unified API to all the biodiversity sources that we provide. The obvious advantage is that a user can interact with a common API and not worry about the nuances in syntax that differ between packages. As more data sources come online, users can access even more data without significant changes to their code. However, it is important to note that spocc will never replicate the full functionality that exists within specific packages. Therefore users with a strong interest in one of the specific data sources listed below would benefit from familiarising themselves with the inner working of the appropriate packages.

Data Sources

spocc currently interfaces with six major biodiversity repositories

- 1. Global Biodiversity Information Facility (rgbif) GBIF is a government funded open data repository with several partner organizations with the express goal of providing access to data on Earth's biodiversity. The data are made available by a network of member nodes, coordinating information from various participant organizations and government agencies.
- 2. Berkeley Ecoengine (ecoengine) The ecoengine is an open API built by the Berkeley Initiative for Global Change Biology. The repository provides access to over 3 million specimens from various Berkeley natural history museums. These data span more than a century and provide access to georeferenced specimens, species checklists, photographs, vegetation surveys and resurveys and a variety of measurements from environmental sensors located at reserves across University of California's natural reserve system.
- iNaturalist (rinat) iNaturalist provides access to crowd sourced citizen science data on species observations.
- 4. VertNet (rvertnet) Similar to rgbif, ecoengine, and rbison (see below), VertNet provides access to more than 80 million vertebrate records spanning a large number of institutions and museums primarly covering four major disciplines (mammology, herpetology, ornithology, and icthyology). Note that we don't currenlty support VertNet data in this package, but we should soon
- 5. Biodiversity Information Serving Our Nation (rbison) Built by the US Geological Survey's core science analytic team, BISON is a portal that provides access to species occurrence data from several participating institutions.
- 6. eBird (rebird) ebird is a database developed and maintained by the Cornell Lab of Ornithology and the National Audubon Society. It provides real-time access to checklist data, data on bird abundance and distribution, and community reports from birders.
- 7. AntWeb (AntWeb) AntWeb is the world's largest online database of images, specimen records, and natural history information on ants. It is community driven and open to contribution from anyone with specimen records, natural history comments, or images.

Note: It's important to keep in mind that several data providers interface with many of the above mentioned repositories. This means that occurrence data obtained from BISON may be duplicates of data that are also

available through GBIF. We do not have a way to resolve these duplicates or overlaps at this time but it is an issue we are hoping to resolve in future versions of the package. See ?spocc_duplicates, after installation, for more.

Data retrieval

The most significant function in spoce is the occ (short for occurrence) function. occ takes a query, often a species name, and searches across all data sources specified in the from argument. For example, one can search for all occurrences of Sharp-shinned Hawks (Accipiter striatus) from the GBIF database with the following R call.

```
library(spocc)
df <- occ(query = 'Accipiter striatus', from = 'gbif')
df

## Summary of results - occurrences found for:
## gbif : 25 records across 1 species
## bison : 0 records across 1 species
## inat : 0 records across 1 species
## ebird : 0 records across 1 species
## ecoengine : 0 records across 1 species
## antweb : 0 records across 1 species</pre>
```

head(df\$gbif\$data[[1]])

```
name
                             key decimalLatitude decimalLongitude prov
## 1 Accipiter striatus 8.91e+08
                                           43.13
                                                           -72.53 gbif
## 2 Accipiter striatus 8.91e+08
                                           32.86
                                                           -97.20 gbif
## 3 Accipiter striatus 8.91e+08
                                                           -71.73 gbif
                                           18.27
                                                           -97.65 gbif
## 4 Accipiter striatus 8.91e+08
                                           30.16
## 5 Accipiter striatus 8.91e+08
                                           37.49
                                                          -122.44 gbif
## 6 Accipiter striatus 8.91e+08
                                           43.63
                                                           -73.07 gbif
```

The data returned are part of a S3 class called occdat. This class has slots for the five data sources described above. One can easily switch the source by changing the from parameter in the function call above.

Within each data source is the set of species queried. In the above example, we only asked for occurrence data for one species, but we could have asked for any number. Let's say we asked for data for two species: *Accipiter striatus*, and *Pinus contorta*. Then the structure of the response would be

If you only request data from gbif, like from = 'gbif', then the other four source slots are present in the response object, but have no data.

You can quickly get just the data by indexing to the data element, like

head(df\$gbif\$data\$Accipiter_striatus)

```
key decimalLatitude decimalLongitude prov
##
                   name
## 1 Accipiter striatus 8.91e+08
                                         43.13
                                                          -72.53 gbif
## 2 Accipiter striatus 8.91e+08
                                          32.86
                                                          -97.20 gbif
## 3 Accipiter striatus 8.91e+08
                                         18.27
                                                          -71.73 gbif
## 4 Accipiter striatus 8.91e+08
                                          30.16
                                                          -97.65 gbif
                                                         -122.44 gbif
## 5 Accipiter striatus 8.91e+08
                                          37.49
## 6 Accipiter striatus 8.91e+08
                                          43.63
                                                          -73.07 gbif
```

When you get data from multiple providers, the fields returned are slightly different, e.g.:

```
df <- occ(query = 'Accipiter striatus', from = c('gbif', 'ecoengine'))
head(df$gbif$data$Accipiter_striatus)</pre>
```

```
##
                            key decimalLatitude decimalLongitude prov
                  name
                                                          -72.53 gbif
## 1 Accipiter striatus 8.91e+08
                                          43.13
                                          32.86
## 2 Accipiter striatus 8.91e+08
                                                         -97.20 gbif
## 3 Accipiter striatus 8.91e+08
                                          18.27
                                                          -71.73 gbif
## 4 Accipiter striatus 8.91e+08
                                          30.16
                                                          -97.65 gbif
## 5 Accipiter striatus 8.91e+08
                                          37.49
                                                         -122.44 gbif
## 6 Accipiter striatus 8.91e+08
                                          43.63
                                                          -73.07 gbif
```

head(df\$ecoengine\$data\$Accipiter_striatus)

```
## longitude latitude type

## 1 -81.52 41.08 Feature

## 2 -135.27 57.53 Feature

## 3 -134.29 57.71 Feature

## 4 -134.29 57.71 Feature

## 5 -148.24 60.55 Feature

## 6 -145.82 61.11 Feature

## 1 https://ecoengine.berkeley.edu/api/observations/MVZ%3ABird%3A127226/
```

```
## 2
        https://ecoengine.berkeley.edu/api/observations/MVZ%3ABird%3A149/
## 3
       https://ecoengine.berkeley.edu/api/observations/MVZ%3ABird%3A268/
## 4
        https://ecoengine.berkeley.edu/api/observations/MVZ%3ABird%3A269/
       https://ecoengine.berkeley.edu/api/observations/MVZ%3ABird%3A1143/
## 5
## 6
       https://ecoengine.berkeley.edu/api/observations/MVZ%3ABird%3A1145/
     observation type
                                                     country state province
##
                                          name
## 1
             specimen Accipiter striatus velox United States
                                                                        Ohio
             specimen Accipiter striatus velox United States
## 2
                                                                      Alaska
## 3
             specimen Accipiter striatus velox United States
                                                                      Alaska
## 4
             specimen Accipiter striatus velox United States
                                                                      Alaska
## 5
             specimen Accipiter striatus velox United States
                                                                      Alaska
## 6
             specimen Accipiter striatus velox United States
                                                                      Alaska
##
                  end_date
     begin_date
                                                                   source
## 1 1867-04-26 1867-04-26 https://ecoengine.berkeley.edu/api/sources/1/
## 2 1907-08-26 1907-08-26 https://ecoengine.berkeley.edu/api/sources/1/
## 3 1907-05-25 1907-05-25 https://ecoengine.berkeley.edu/api/sources/1/
## 4 1907-05-24 1907-05-24 https://ecoengine.berkeley.edu/api/sources/1/
## 5 1908-08-18 1908-08-18 https://ecoengine.berkeley.edu/api/sources/1/
## 6 1908-09-03 1908-09-03 https://ecoengine.berkeley.edu/api/sources/1/
                                        remote resource
## 1 http://arctos.database.museum/guid/MVZ:Bird:127226 ecoengine
        http://arctos.database.museum/guid/MVZ:Bird:149 ecoengine
## 3
        http://arctos.database.museum/guid/MVZ:Bird:268 ecoengine
## 4
        http://arctos.database.museum/guid/MVZ:Bird:269 ecoengine
## 5
       http://arctos.database.museum/guid/MVZ:Bird:1143 ecoengine
## 6
       http://arctos.database.museum/guid/MVZ:Bird:1145 ecoengine
```

We provide a function occ2df that pulls out a few key columns needed for making maps:

head(occ2df(df))

```
##
                   name longitude latitude prov
## 1 Accipiter striatus
                           -72.53
                                      43.13 gbif
                           -97.20
                                      32.86 gbif
## 2 Accipiter striatus
## 3 Accipiter striatus
                           -71.73
                                      18.27 gbif
## 4 Accipiter striatus
                           -97.65
                                      30.16 gbif
## 5 Accipiter striatus
                          -122.44
                                      37.49 gbif
## 6 Accipiter striatus
                           -73.07
                                      43.63 gbif
```

Fix names

One problem you often run in to is that there can be various names for the same taxon in any one source. For example:

```
df <- occ(query='Pinus contorta', from=c('gbif','inat'), limit = 50)
head(df$gbif$data$Pinus_contorta[,1:2])</pre>
```

```
## 5 Pinus contorta 899974349
## 6 Pinus contorta 910497271
```

head(df\$inat\$data\$Pinus_contorta[,1:2])

This is fine, but when trying to make a map in which points are colored for each taxon, you can have many colors for a single taxon, where instead one color per taxon is more appropriate. There is a function in spocc called fixnames, which has a few options in which you can take the shortest names (usually just the plain binomials like *Homo sapiens*), or the original name queried, or a vector of names supplied by the user.

```
df <- fixnames(df, how = 'shortest')
head(df$gbif$data$Pinus_contorta[,1:2])</pre>
```

head(df\$inat\$data\$Pinus_contorta[,1:2])

```
df_comb <- occ2df(df)
head(df_comb); tail(df_comb)</pre>
```

```
name longitude latitude prov
## 1 Pinus contorta
                        -122.8
                                  48.14 gbif
## 2 Pinus contorta
                        -120.0
                                  38.82 gbif
## 3 Pinus contorta
                        -120.3
                                  39.34 gbif
## 4 Pinus contorta
                        -120.2
                                  39.20 gbif
## 5 Pinus contorta
                        -122.3
                                  47.66 gbif
## 6 Pinus contorta
                        -120.2
                                  39.32 gbif
```

```
##
       name longitude latitude prov
                           39.42 inat
## 95
                -120.6
                -120.7
## 96
                           39.42 inat
## 97
                -121.6
                           46.90 inat
## 98
                -123.0
                           48.47 inat
## 99
                -121.8
                           46.75 inat
## 100
                -120.2
                           39.43 inat
```

Visualization routines

Interactive maps

Leaflet.js

Leaflet JS is an open source mapping library that can leverage various layers from multiple sources. Using the leafletR library, it's possible to generate a local geoJSON file and a html file of species distribution maps. The folder can easily be moved to a web server and served widely without any additional coding.

It's also possible to render similar maps with Mapbox by committing just the geoJSON file to GitHub or posting it as a gist on GitHub. All the remaining fields will become part of a table inside a tooltip, providing a extremely quick and easy way to serve up interactive maps. This is especially useful when users do not have their own web hosting options.

Here is an example of making a leaflet map:

```
spp <- c('Danaus plexippus','Accipiter striatus','Pinus contorta')
dat <- occ(query = spp, from = 'gbif', gbifopts = list(hasCoordinate = TRUE))
data <- occ2df(dat)
mapleaflet(data = data, dest = ".")</pre>
```



Geojson map as a Github gist

You can also create interactive maps via the mapgist function. You have to have a Github account to use this function. Github accounts are free though, and great for versioning and collaborating on code or papers. When you run the mapgist function it will ask for your Github username and password. You can alternatively store those in your .Rprofile file by adding entries for username (options(github.username = 'username')) and password (options(github.password = 'password')).

```
spp <- c('Danaus plexippus','Accipiter striatus','Pinus contorta')
dat <- occ(query = spp, from = 'gbif', gbifopts = list(hasCoordinate = TRUE))
dat <- fixnames(dat)</pre>
```

```
dat <- occ2df(dat)
mapgist(data = dat, color = c("#976AAE", "#6B944D", "#BD5945"))</pre>
```

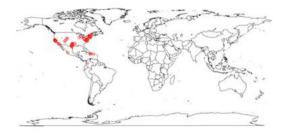


Static maps

base plots

Base plots, or the built in plotting facility in R accessed via plot(), is quite fast, but not easy or efficient to use, but are good for a quick glance at some data.

```
spnames <- c('Accipiter striatus', 'Setophaga caerulescens', 'Spinus tristis')
out <- occ(query = spnames, from = 'gbif', gbifopts = list(hasCoordinate = TRUE))
plot(out, cex = 1, pch = 10)</pre>
```



ggplot2

ggplot2 is a powerful package for making visualizations in R. Read more about it here. We created a simple wrapper function mapggplot to make a ggplot2 map from occurrence data using the ggmap package, which is built on top of ggplot2. Here's an example:

```
ecoengine_data <- occ(query = 'Lynx rufus californicus', from = 'ecoengine')
mapggplot(ecoengine_data)</pre>
```

Upcoming features

- As soon as we have an updated rvertnet package, we'll add the ability to query VertNet data from spocc.
- We will add rCharts as an official import once the package is on CRAN (Eta end of March)



- We're helping on a new package rMaps to make interactive maps using various Javascript mapping libraries, which will give access to a variety of awesome interactive maps. We will integrate rMaps once it's on CRAN.
- We'll add a function to make interactive maps using RStudio's Shiny in a future version.