# auk

extracting & processing
## eBird Data in R

Matt Strimas-Mackey

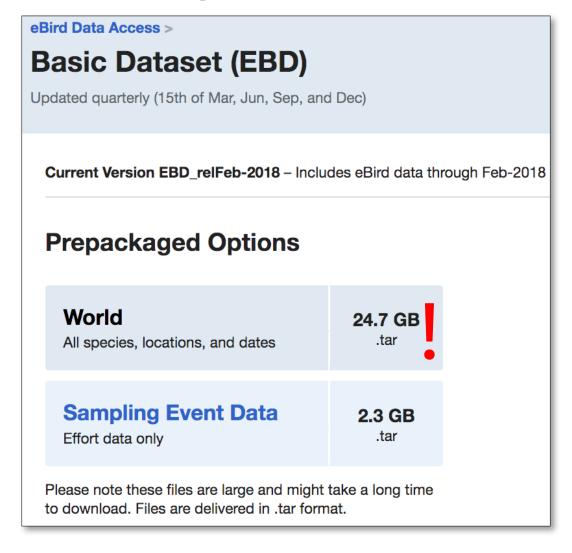*Cornell Lab of Ornithology • strimas.com • github.com/mstrimas*

# 500 million observations

**253**

# countries

# 10,374
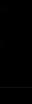## species

# SHOW ME THE DATA!!!

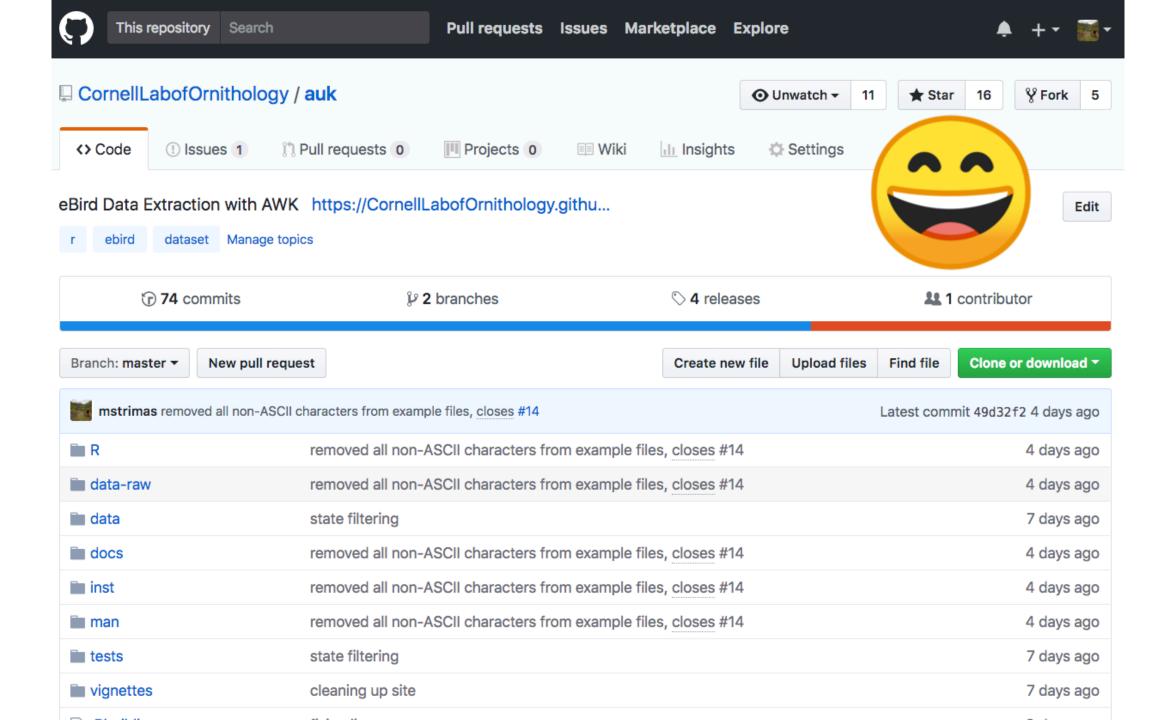we need to make this

# File

smaller

```
ebd_relFeb-2018 — -bash — 116×28
```

```
[ag-clo-mes335:ebd_relFeb-2018 mes335$ cd /Users/mes335/data/ebird/ebd_relFeb-2018
[ag-clo-mes335:ebd_relFeb-2018 mes335$ ls
awk                                  ebd_relFeb-2018_subset.txt
ebd_relFeb-2018_SG.txt               ebd_sampling_relFeb-2018_SG.txt
ebd_relFeb-2018_clean.txt            ebd_sampling_relFeb-2018_clean.txt
ebd_relFeb-2018_rollup.txt           ebird_country-state.txt
[ag-clo-mes335:ebd_relFeb-2018 mes335$ head -20 ebd_relFeb-2018_clean.txt | awk -F $'\t' '{print $25,$26,$27,$5,$9'}
LATITUDE LONGITUDE OBSERVATION DATE COMMON NAME OBSERVATION COUNT
35.1085767 -120.6265268 2008-11-05 Belted Kingfisher 1
35.1085767 -120.6265268 2008-11-05 Wilson's Snipe 1
35.1085767 -120.6265268 2008-11-05 Ruby-crowned Kinglet 15
35.1085767 -120.6265268 2008-11-05 California Gull 35
35.1085767 -120.6265268 2008-11-05 Mallard 40
35.1085767 -120.6265268 2008-11-05 Bewick's Wren 1
8.52111 -83.39991 2004-06-18 Rufous Piha X
8.52111 -83.39991 2004-06-18 Thick-billed Euphonia X
46.3024735 -91.0284233 2008-08-16 Chestnut-sided Warbler 4
46.3024735 -91.0284233 2008-08-16 American Robin 9
33.7968382 -117.7532673 1971-01-02 Spotted Towhee 4
33.7968382 -117.7532673 1971-01-02 Red-tailed Hawk 5
33.7968382 -117.7532673 1971-01-02 Anna's Hummingbird 11
33.7968382 -117.7532673 1971-01-02 Bewick's Wren 5
37.1434192 -114.0232372 2008-10-25 Spotted Towhee 6
43.1812545 -88.0415756 2008-08-06 Mallard 5
43.1812545 -88.0415756 2008-08-06 Red-winged Blackbird 3
41.059032 -112.238045 2006-08-13 Black-necked Stilt 1
41.059032 -112.238045 2006-08-13 California Gull 1
ag-clo-mes335:ebd_relFeb-2018 mes335$ 
```

Terminal 😑

AWK 😟

```
awk-example.txt
1
2   BEGIN {
3       FS = OFS = "  "
4           split("Cyanocitta cristata  Perisoreus canadensis", speciesValues, "  ")
5           for (i in speciesValues) species[speciesValues[i]] = 1
6           split("CA US", countryValues, "  ")
7           for (i in countryValues) countries[countryValues[i]] = 1
8   }
9   {
10      keep = 1
11      if (keep == 1 && ($6 in species)) {
12          keep = 1
13      } else {
14          keep = 0
15      }
16      if (keep == 1 && ($14 in countries)) {
17          keep = 1
18      } else {
19          keep = 0
20      }
21      if (keep == 1 && ($26 >= -100 && $26 <= -80 && $25 >= 37 && $25 <= 52)) {
22          keep = 1
23      } else {
24          keep = 0
25      }
26      if (keep == 1 && ($27 >= "2012-01-01" && $27 <= "2012-12-31")) {
27          keep = 1
28      } else {
29          keep = 0
30      }
31      if (keep == 1 && ($28 >= "06:00" && $28 <= "09:00")) {
32          keep = 1
```

CornellLabofOrnithology / auk

Unwatch ▾ 11    ★ Star 16    Fork 5

<> Code    ⓘ Issues 1    Pull requests 0    Projects 0    Wiki    Insights    Settings

eBird Data Extraction with AWK    https://CornellLabofOrnithology.githu…    Edit

r    ebird    dataset    Manage topics

⊙ 74 commits    2 branches    ⬙ 4 releases    1 contributor

Branch: master ▾    New pull request    Create new file    Upload files    Find file    Clone or download ▾

mstrimas removed all non-ASCII characters from example files, closes #14    Latest commit 49d32f2 4 days ago

| | | |
|---|---|---|
| R | removed all non-ASCII characters from example files, closes #14 | 4 days ago |
| data-raw | removed all non-ASCII characters from example files, closes #14 | 4 days ago |
| data | state filtering | 7 days ago |
| docs | removed all non-ASCII characters from example files, closes #14 | 4 days ago |
| inst | removed all non-ASCII characters from example files, closes #14 | 4 days ago |
| man | removed all non-ASCII characters from example files, closes #14 | 4 days ago |
| tests | state filtering | 7 days ago |
| vignettes | cleaning up site | 7 days ago |

# auk workflow

1. clean
2. filter
3. import
4. pre-process
5. zero-fill

**1** clean

remove problematic rows
drop free-text columns

**2** **filter**

extract just what you need
reduce file to a manageable size

**3** import

import to a data frame
choose clean names and sensible data types

# pre-process

## taxonomic roll-up

| | |
|---|---|
| 25 | Common Goldeneye |
| 2000 | duck sp. |
| 40 | Mallard (Northern) |
| 20 | Mallard/American Black Duck |
| 6 | Glaucous-winged Gull |
| 4 | Western x Glaucous-winged Gull (hybrid) |
| 55 | Rock Pigeon (Feral Pigeon) |
| 25 | Yellow-rumped Warbler (Myrtle) |
| 90 | Yellow-rumped Warbler (Audubon's) |

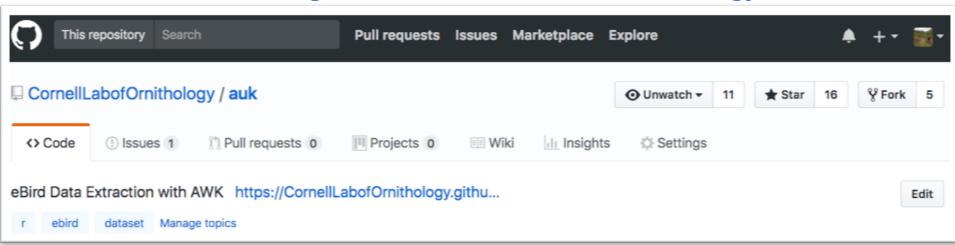## 5 zero-fill

complete checklists give implicit zeros
use sampling event data to make zeros explicit

Read the vignette: CornellLabofOrnithology.github.io/auk



File an issue: github.com/CornellLabofOrnithology/auk



Send me an email: mes335@cornell.edu