

Alakazam: Analysis of clonal abundance and diversity

Jason Anthony Vander Heiden

2018-09-13

Contents

Example data	1
Generate a clonal abundance curve	2
Generate a diversity curve	3
Test diversity at a fixed diversity order	5

The clonal diversity of the repertoire can be analyzed using the general form of the diversity index, as proposed by Hill in:

Hill, M. Diversity and evenness: a unifying notation and its consequences.
Ecology 54, 427-432 (1973).

Coupled with resampling strategies to correct for variations in sequencing depth, as well as inference of complete clonal abundance distributions as described in:

Chao A, et al. Rarefaction and extrapolation with Hill numbers:
A framework for sampling and estimation in species diversity studies.
Ecol Monogr. 2014 84:45-67.

Chao A, et al. Unveiling the species-rank abundance distribution by
generalizing the Good-Turing sample coverage theory.
Ecology. 2015 96, 11891201.

This package provides methods for the inference of a complete clonal abundance distribution, using the `estimateAbundance` function, along with two approaches to assess diversity of these distributions:

1. Generation of a smooth diversity (D) curve over a range of diversity orders (q) using `rarefyDiversity`.
2. A significance test of the diversity (D) at a fixed diversity order (q) using `testDiversity`.

Example data

A small example Change-O database, `ExampleDb`, is included in the `alakazam` package. Diversity calculation requires the `CLONE` field (column) to be present in the Change-O file, as well as an additional grouping column. In this example we will use the grouping columns `SAMPLE` and `ISOTYPE`.

```
# Load required packages
```

```
library(alakazam)
```

```
# Load example data
```

```
data(ExampleDb)
```

Generate a clonal abundance curve

A simple table of the observed clonal abundance counts and frequencies may be generated using the `countClones` function either without copy numbers, where the size of each clone is determined by the number of sequence members:

```
# Partitions the data based on the SAMPLE column
clones <- countClones(ExampleDb, groups="SAMPLE")
head(clones, 5)

## # A tibble: 5 x 4
## # Groups:   SAMPLE [1]
##   SAMPLE CLONE SEQ_COUNT SEQ_FREQ
##   <chr>  <chr>      <int>   <dbl>
## 1 +7d    3128         100    0.1
## 2 +7d    3100          50    0.05
## 3 +7d    3141          44    0.044
## 4 +7d    3177          30    0.03
## 5 +7d    3170          28    0.028
```

You may also specify a column containing the abundance count of each sequence (usually copy numbers), that will including weighting of each clone size by the corresponding abundance count. Furthermore, multiple grouping columns may be specified such that `SEQ_FREQ` (unweighted clone size as a fraction of total sequences in the group) and `COPY_FREQ` (weighted fraction) are normalized to within multiple group data partitions.

```
# Partitions the data based on both the SAMPLE and ISOTYPE columns
# Weights the clone sizes by the DUPCOUNT column
clones <- countClones(ExampleDb, groups=c("SAMPLE", "ISOTYPE"), copy="DUPCOUNT")
head(clones, 5)

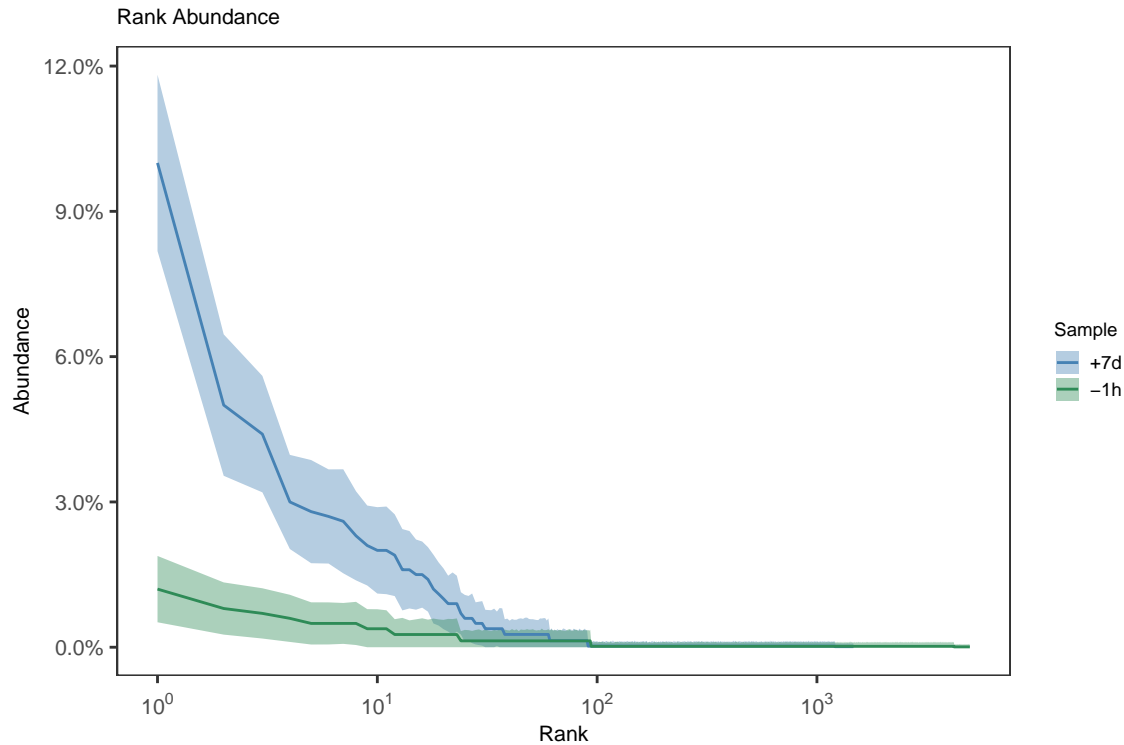
## # A tibble: 5 x 7
## # Groups:   SAMPLE, ISOTYPE [2]
##   SAMPLE ISOTYPE CLONE SEQ_COUNT COPY_COUNT SEQ_FREQ COPY_FREQ
##   <chr>  <chr>  <chr>    <int>    <int>    <dbl>   <dbl>
## 1 +7d    IgA     3128     88      651    0.331   0.497
## 2 +7d    IgG     3100     49      279    0.0928  0.173
## 3 +7d    IgA     3141     44      240    0.165   0.183
## 4 +7d    IgG     3192     19      141    0.0360  0.0874
## 5 +7d    IgG     3177     29      130    0.0549  0.0806
```

While `countClones` will report observed abundances, it will not correct the distribution nor provide confidence intervals. A complete clonal abundance distribution may be inferred using the `estimateAbundance` function with confidence intervals derived via bootstrapping. This output may be visualized using the `plotAbundanceCurve` function.

```
# Partitions the data on the SAMPLE column
# Calculates a 95% confidence interval via 200 bootstrap realizations
clones <- estimateAbundance(ExampleDb, group="SAMPLE", ci=0.95, nboot=200)

# Plots a rank abundance curve of the relative clonal abundances
```

```
sample_colors <- c("-1h"="seagreen", "+7d"="steelblue")
plot(clones, colors=sample_colors, legend_title="Sample")
```



Generate a diversity curve

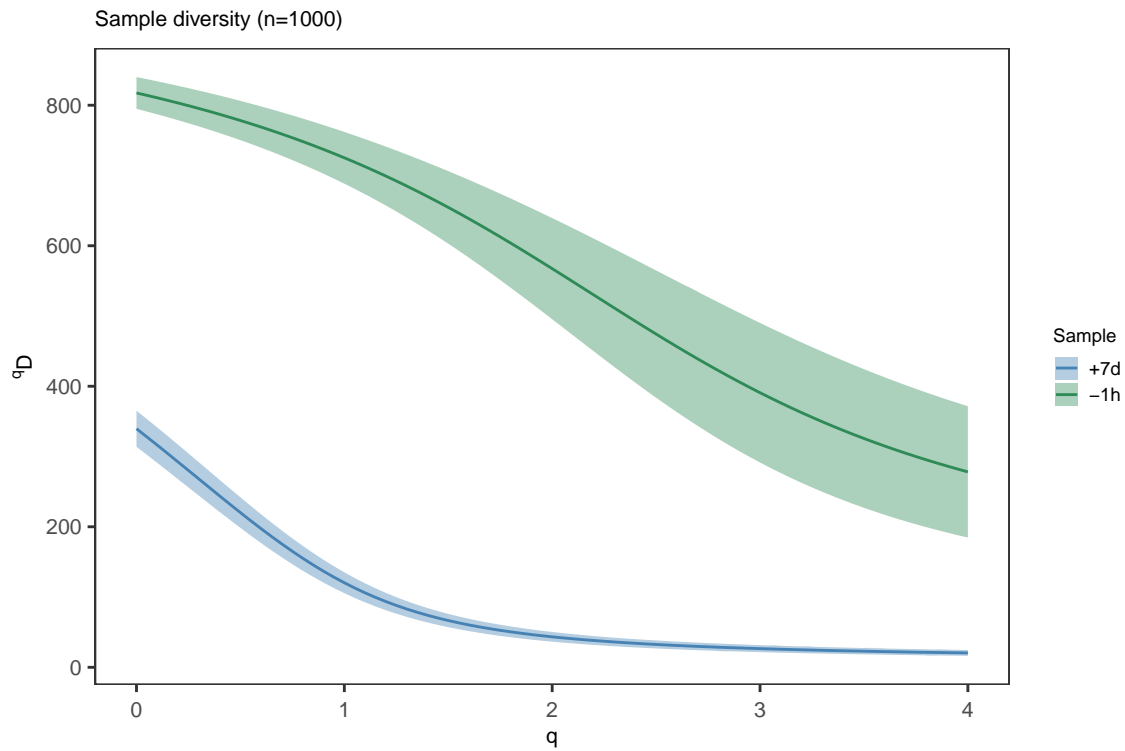
The function `rarefyDiversity` performs uniform resampling of the input sequences and recalculates the clone size distribution, and diversity, with each resampling realization. Diversity (D) is calculated over a range of diversity orders (q) to generate a smooth curve.

```
# Compare diversity curve across values in the "SAMPLE" column
# q ranges from 0 (min_q=0) to 4 (max_q=4) in 0.05 increments (step_q=0.05)
# A 95% confidence interval will be calculated (ci=0.95)
# 200 resampling realizations are performed (nboot=200)
sample_div <- rarefyDiversity(ExampleDb, "SAMPLE", min_q=0, max_q=4, step_q=0.05,
                             ci=0.95, nboot=200)

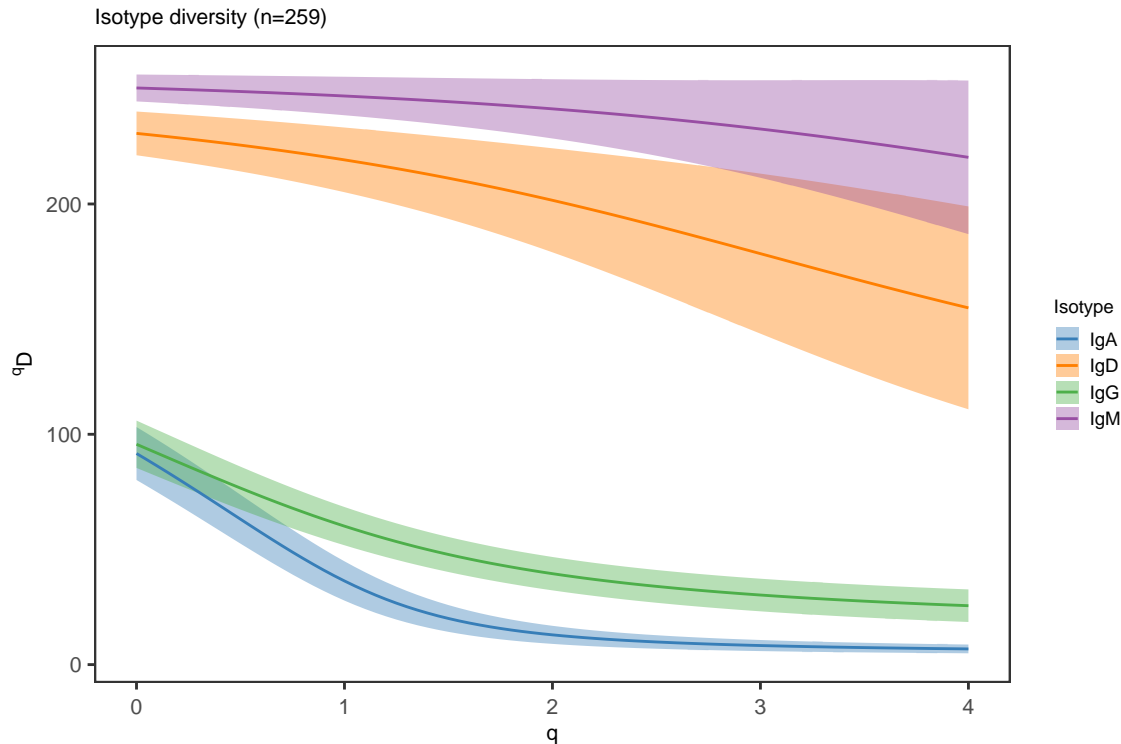
# Compare diversity curve across values in the "ISOTYPE" column
# Analyse is restricted to ISOTYPE values with at least 30 sequences by min_n=30
# Excluded groups are indicated by a warning message
isotype_div <- rarefyDiversity(ExampleDb, "ISOTYPE", min_n=30, min_q=0, max_q=4,
                              step_q=0.05, ci=0.95, nboot=200)

# Plot a log-log (log_q=TRUE, log_d=TRUE) plot of sample diversity
# Indicate number of sequences resampled from each group in the title
sample_main <- paste0("Sample diversity (n=", sample_div@n, ")")
sample_colors <- c("-1h"="seagreen", "+7d"="steelblue")
```

```
plot(sample_div, colors=sample_colors, main_title=sample_main,
      legend_title="Sample")
```



```
# Plot isotype diversity using default set of Ig isotype colors
isotype_main <- paste0("Isotype diversity (n=", isotype_div@n, ")")
plot(isotype_div, colors=IG_COLORS, main_title=isotype_main,
      legend_title="Isotype")
```



Test diversity at a fixed diversity order

The function `testDiversity` performs resampling and diversity calculation in the same manner as `rarefyDiversity`, but only for a single diversity order. Significance testing across groups is performed using the delta of the bootstrap distributions between groups.

```
# Test diversity at q=2 (equivalent to Simpson's index) across values in the "SAMPLE" column
# 200 bootstrap realizations are performed (nboot=200)
sample_test <- testDiversity(ExampleDb, 2, "SAMPLE", nboot=200)
```

```
# Print p-value table
print(sample_test)
```

```
##          test DELTA_MEAN DELTA_SD PVALUE
## 1 +7d != -1h  522.2804 37.80454      0
```

```
# Test diversity across values in the "ISOTYPE" column
# Analyse is restricted to ISOTYPE values with at least 30 sequences by min_n=30
# Excluded groups are indicated by a warning message
isotype_test <- testDiversity(ExampleDb, 2, "ISOTYPE", min_n=30, nboot=200)
```

```
# Print p-value table
print(isotype_test)
```

```
##          test DELTA_MEAN DELTA_SD PVALUE
## 1 IgA != IgD  189.10274 11.194808      0
```

```
## 2 IgA != IgG 26.88033 4.627724 0
## 3 IgA != IgM 228.89622 6.651614 0
## 4 IgD != IgG 162.22241 11.629304 0
## 5 IgD != IgM 39.79348 13.018352 0
## 6 IgG != IgM 202.01589 8.095055 0
```

```
# Plot the mean and standard deviations from the test
plot(isotype_test, colors=IG_COLORS, main_title=isotype_main,
      legend_title="Isotype")
```

