

Alakazam: Analysis of clonal abundance and diversity

Jason Anthony Vander Heiden

2016-07-22

Contents

Example data	1
Generate a clonal abundance curve	2
Generate a diversity curve	4
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)
```

```
## Source: local data frame [5 x 4]
## Groups: SAMPLE [1]
##
##   SAMPLE CLONE SEQ_COUNT SEQ_FREQ
##   <chr> <chr>    <int>    <dbl>
## 1    +7d  3128      100    0.100
## 2    +7d  3100       50    0.050
## 3    +7d  3141       44    0.044
## 4    +7d  3177       30    0.030
## 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)
```

```
## Source: local data frame [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 0.33082707 0.49732620
## 2    +7d    IgG  3100      49      279 0.09280303 0.17296962
## 3    +7d    IgA  3141      44      240 0.16541353 0.18334607
## 4    +7d    IgG  3192      19      141 0.03598485 0.08741476
## 5    +7d    IgG  3177      29      130 0.05492424 0.08059516
```

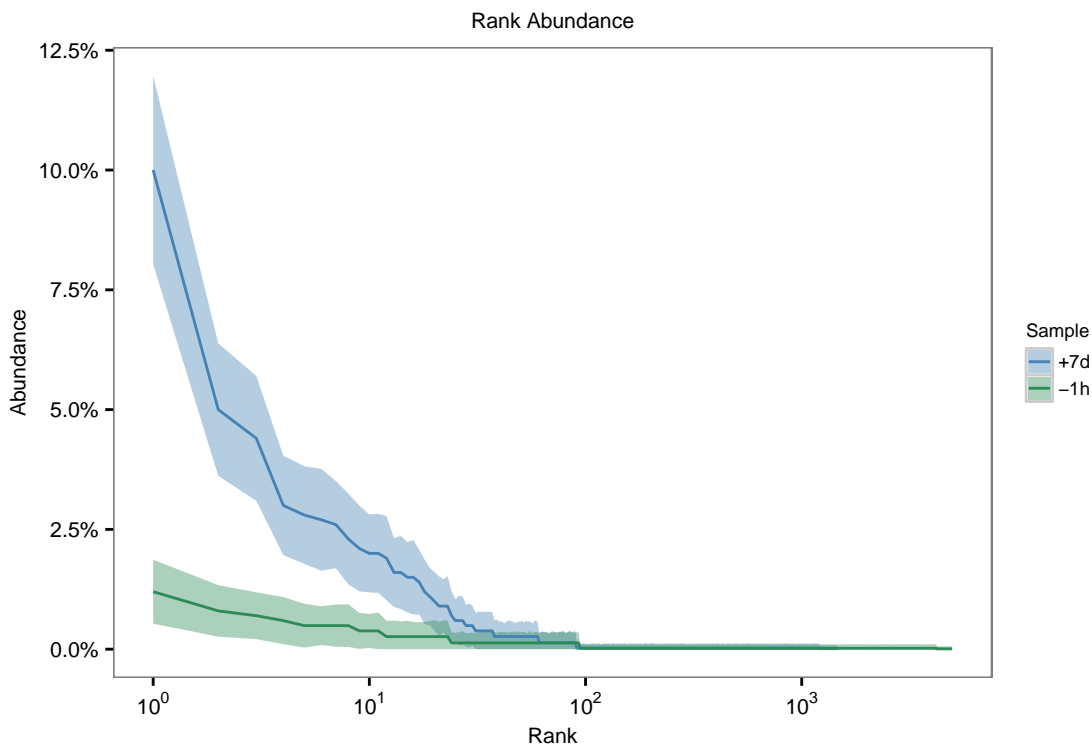
While `countClones` will report observed abundances, it will not correct the distribution nor provide confidence intervals. A complete clonal abundance distribution may be inferred using the `estimateAbundance` function with confidence intervals derived via bootstrapping. This output may be visualized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

```
head(clones, 5)
```

```
## # A tibble: 5 x 6
##   GROUP CLONE      P      LOWER      UPPER  RANK
##   <chr> <chr> <dbl>    <dbl>    <dbl> <int>
## 1   +7d  3128 0.100 0.08034570 0.11965430     1
## 2   +7d  3100 0.050 0.03618494 0.06381506     2
## 3   +7d  3141 0.044 0.03096411 0.05703589     3
## 4   +7d  3177 0.030 0.01963274 0.04036726     4
## 5   +7d  3170 0.028 0.01781346 0.03818654     5
```

```
# Plots a rank abundance curve of the relative clonal abundances
sample_colors <- c("-1h"="seagreen", "+7d"="steelblue")
plotAbundance(clones, colors=sample_colors, legend_title="Sample")
```



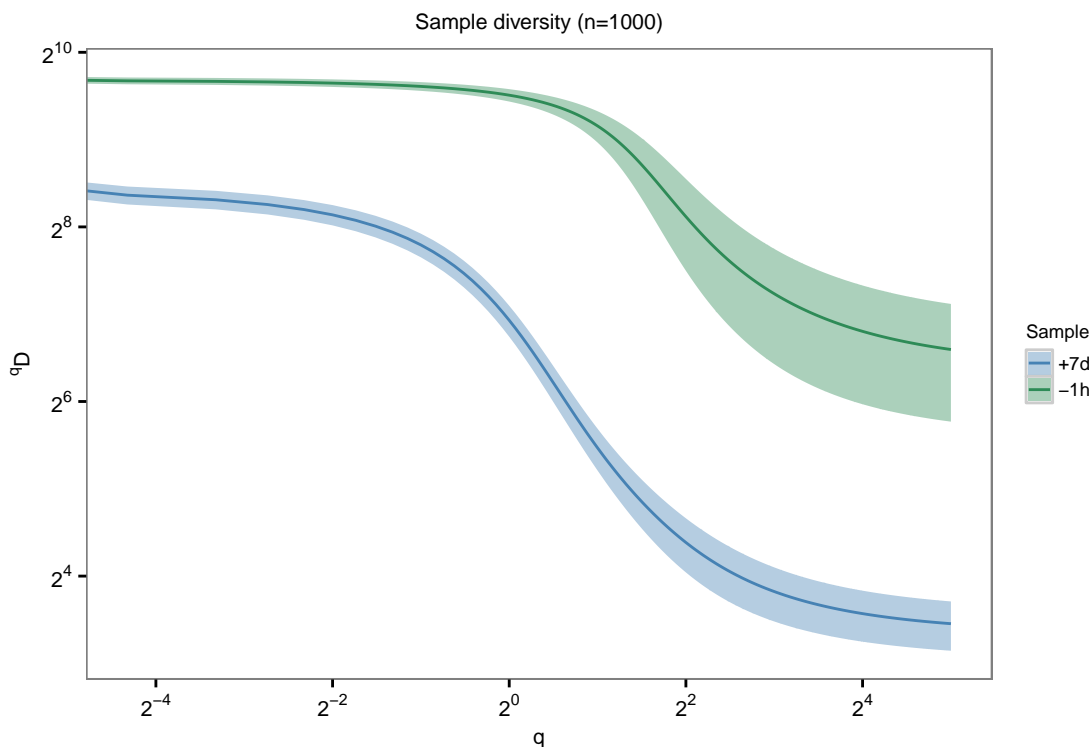
Generate a diversity curve

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.

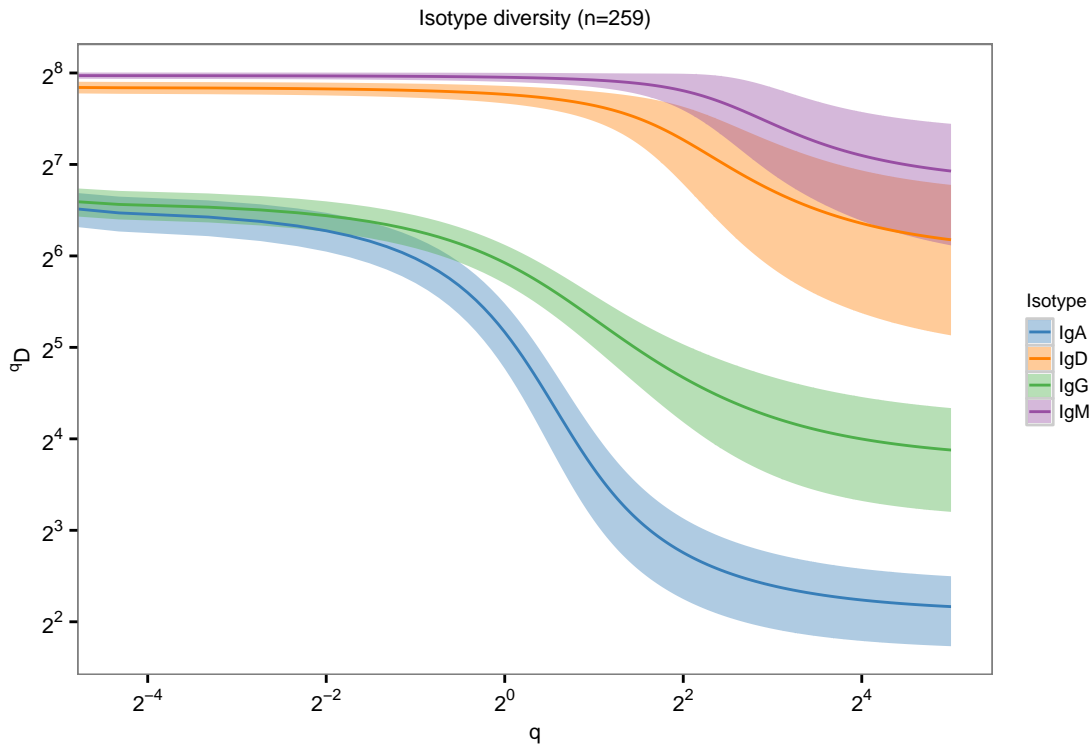
```
# Compare diversity curve across values in the "SAMPLE" column
# q ranges from 0 (min_q=0) to 32 (max_q=32) in 0.05 increments (step_q=0.05)
# A 95% confidence interval will be calculated (ci=0.95)
# 2000 resampling realizations are performed (nboot=200)
sample_div <- rarefyDiversity(ExampleDb, "SAMPLE", min_q=0, max_q=32, 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=32,
                              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")
plotDiversityCurve(sample_div, colors=sample_colors, main_title=sample_main,
                  legend_title="Sample", log_q=TRUE, log_d=TRUE)
```



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



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=0 (species richness) across values in the "SAMPLE" column
# 2000 bootstrap realizations are performed (nboot=200)
```

```
sample_test <- testDiversity(ExampleDb, 0, "SAMPLE", nboot=200)
```

```
# Print p-value table
```

```
print(sample_test)
```

```
##          test DELTA_MEAN DELTA_SD PVALUE
## 1 +7d != -1h    475.015 17.35687      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  189.12401  11.283572      0
## 2 IgA != IgG   26.17042   5.009623      0
## 3 IgA != IgM  228.78501   6.771065      0
## 4 IgD != IgG  162.95359  11.881203      0
## 5 IgD != IgM   39.66100  12.602362      0
## 6 IgG != IgM  202.61459   8.091725      0
```