

Phanerozoic-scale global marine biodiversity analysis with the R package divDyn v0.6

Adam T. Kocsis, John Alroy, Carl J. Reddin, Wolfgang Kiessling

2018-09-20

1. Introduction

In this text we describe how to use the ‘divDyn’ R package for Phanerozoic scale global biodiversity analysis with data acquired from the Paleobiology Database (PaleoDB, <http://www.paleobiodb.org/>, <http://fossilworks.org>) until the final results. The aim of this vignette is to describe a fully reproducible workflow. This text does not explain all the capabilities of the package. For a full treatment, please refer to the more detailed vignette ‘Handout to the R package ‘divDyn’ v0.6.0 for diversity dynamics from fossil occurrence data’.

1.1. Necessary files and installation

Although the source code has been submitted to the CRAN servers, it might take some time until the archive is completely processed and you can run

```
install.packages("divDyn")
```

to install divDyn. If this fails, check the package GitHub repository

<http://www.github.com/adamkocsis/divDyn>

for further instructions.

The files necessary to replicate the analyses presented here are deposited in another dedicated GitHub repository

<http://www.github.com/adamkocsis/ddPhanero>

and are referenced directly from the vignette. We intend to update this material for future reference. The analytical script presented here is also available as a regular .R file (scripts/0.3/2018-09-14_marineAnimals_ddPhanero.R).

2. Getting the data

The occurrence data used in the study were retrieved from the Paleobiology Database. We used all occurrences from the Ediacaran to Holocene interval, with several additional variables (lithology, environment) for further analyses.

2.1. Our downloaded data

To ensure the exact reproducibility of the material presented here, the raw dataset that was used to calculate the results (downloaded on September 14, 2018) is available for download and experimentation. If you do not wish to carry out a PaleoDB download yourself, you can use our file.

```
load(url(
  "https://github.com/adamkocsis/ddPhanero/raw/master/data/PaleoDB/2018-09-14_paleoDB.RData"
))
```

This particular download has

```
nrow(dat)
```

```
## [1] 1369299
```

rows. The command above will import a single `data.frame` class object `dat` to the global workspace. The metadata are available from the GitHub repository of the example analysis.

2.2. New download

In case you do want to carry out a new download, the download query is reproducible with the link that we used to access the data. The download will commence automatically if you copy and paste this URL into your browser.

```
https://paleobiodb.org/data1.2/occs/list.csv?interval=Ediacaran,Holocene&
show=class,classext,genus,subgenus,abund,coll,coords,loc,paleoloc,strat,stratext,lith,
env,ref,crmod, timebins,timecompare
```

It will take a considerable amount of time (30+ minutes) to download this set, and the resulting file is very large (ca. 1.4 GB). Nevertheless, it can be stored effectively by omitting some variables and saving it as a binary `.RData` format. You can load the file (saved as “allData_2018-09-14.csv” in this case) using the `read.csv()` function.

```
# assuming that the file is in your current working directory
dat <- read.csv("allData_2018-09-14.csv", header=TRUE, stringsAsFactors=FALSE)
```

Then you can define the required columns. We used the following set so the final file can be retrieved from GitHub without timeout problems.

```
need<-c("collection_no", "collection_name", "accepted_name",
  "accepted_rank", "identified_name", "identified_rank",
  "early_interval", "late_interval", "max_ma", "min_ma",
  "reference_no", "phylum", "class", "order", "family",
  "genus", "lng", "lat", "paleolng", "paleolat", "formation",
  "lithology1", "lithification1", "environment", "created",
  "zone")
dat <- dat[,need]
```

Then you can save everything as an `.RData` file,

```
save(dat, file="allData_2018-09-14.RData")
```

and use `load()` to read the file whenever you need it.

```
load(file="allData_2018-09-14.RData")
```

All additional filtering and transformation will be explained in the following section.

3. Data processing

Processing the raw download consists of steps that can have a huge effect on the results. These transformations include basic taxonomic and environmental filtering, as well as assigning the occurrences to stratigraphic bins

and environmental categories.

3.1. Taxonomic filtering

This dataset is filtered to only comprise marine animal taxa and heterotrophic protists, i.e. the same taxonomic groups listed in Sepkoski's (2002) compendium. As data downloads can be tedious when the taxonomic filters of the PaleoDB are used, we decided to omit taxa procedurally using the available higher-level taxonomic fields: `phylum`, `class`, `order` and `family`.

Filtering started with omission of occurrences that were not identified to the level of genus, which is indicated in the `accepted_rank` field:

```
dat <- dat[dat$accepted_rank %in% c("genus", "species"), ]
```

Omitting these poorly identified fossils decreases the sample size to

```
nrow(dat)
```

```
## [1] 1197589
```

occurrences. The filtering process continues with the selection of phyla that have or had recorded marine species. Phyla that have considerable terrestrial and marine records (chordates and arthropods) are included based on lower taxonomic ranks.

```
# sampled phyla
# levels(factors(dat$phylum))
#A. phyla
marineNoPlant <- c("",
  "Agmata",
  "Annelida",
  "Bilateralomorpha",
  "Brachiopoda",
  "Bryozoa",
  "Calcispongea",
  "Chaetognatha",
  "Cnidaria",
  "Ctenophora",
  "Echinodermata",
  "Entoprocta",
  "Foraminifera",
  "Hemichordata",
  "Hyolitha",
  "Mollusca",
  "Nematoda",
  "Nematomorpha",
  "Nemertina",
  "Onychophora",
  "Petalonamae",
  "Phoronida",
  "Platyhelminthes",
  "Porifera",
  "Problematica",
  "Rhizopodea",
  "Rotifera",
  "Sarcomastigophora",
  "Sipuncula",
```

```

"Uncertain",
"Vetulicolia",
""
)

```

Then a logical vector was defined for each row, suggesting whether the phylum of that occurrence is present or not in the above defined set `marineNoPlant`.

```

# logical vector of rows indicating these
bByPhyla <- dat$phylum %in% marineNoPlant

```

The rest of the data were saved for further filtering. Classes that are likely to be marine animals were extracted from this subset.

```

# noNeed <- dat[!bByPhyla,]
#B. classes
#levels(factor(noNeed$class))
needClass <- c(
  "Acanthodii",
  "Actinopteri",
  "Actinopterygii",
  "Agnatha",
  "Cephalaspidomorphi",
  "Chondrichthyes",
  "Cladistia",
  "Coelacanthimorpha",
  "Conodonts",
  "Galeaspida",
  "Myxini",
  "Osteichthyes",
  "Petromyzontida",
  "Plagiostomi",
  "Pteraspdomorphi",
  "Artiopoda",
  "Branchiopoda",
  "Cephalocarida",
  "Copepoda",
  "Malacostraca",
  "Maxillopoda",
  "Megacheira",
  "Merostomoidea",
  "Ostracoda",
  "Paratrilobita",
  "Pycnogonida",
  "Remipedia",
  "Thylacocephala",
  "Trilobita",
  "Xiphosura"
)
# logical vector of rows indicating occurrences
bNeedClass <- dat$class %in% needClass

```

The mammalian orders Sirenia and Cetacea were also included, as well as pinniped families from the order Carnivora.

```
#C. mammals
# mammals <- dat[dat$class=="Mammalia", ]
# levels(factor(mammals$order))
needMammalOrd <- c("Cetacea", "Sirenia")
bMammalOrder <- dat$order %in% needMammalOrd

# the carnivores
# carnivores <- dat[dat$order=="Carnivora", ]
# levels(factor(carnivores$family))
needFam <- c("Otariidae", "Phocidae", "Desmatophocidae")
bNeedMamFam <- dat$family %in% needFam
```

Some reptile orders were also included:

```
# D. Reptiles
# reptiles <- dat[dat$class=="Reptilia", ]
# levels(factor(reptiles$order))

needReptOrd<-c(
  "Eosauropterygia",
  "Hupehsuchia",
  "Ichthyosauria",
  "Placodontia",
  "Sauropterygia",
  "Thalattosauria"
)
# the logical vector for the total data
bRept <- dat$order %in% needReptOrd
```

Families of sea turtles are also included in the analyzed set.

```
# E. Sea turtles
# turtles <- dat[dat$order=="Testudines", ]
# levels(factor(turtles$family))

needTurtleFam <- c(
  "Cheloniidae",
  "Protostegidae",
  "Dermochelyidae",
  "Dermochelyoidae",
  "Toxochelyidae",
  "Pancheloniidae"
)
bTurtle <- dat$family%in%needTurtleFam
```

And then, we subsetting the data with these multiple filters. In this sense, the logical OR | works as a union operator between the taxonomic groups.

```
dat <- dat[bByPhyla | bNeedClass | bMammalOrder | bNeedMamFam | bRept | bTurtle, ]
```

The entire procedure decreases the number of occurrences to around 900,000.

```
# the number of rows after taxonomic filtering
nrow(dat)
```

```
## [1] 902032
```

After the filtering by taxonomy was finished, we can make sure that potential homonymies will not affect the results by combining the class names and genus names to create individual entries:

```
# resolve the potential homonymy problem
dat$clgen <- paste(dat$class, dat$genus)
```

3.2. Filtering by sedimentary environment

Some of the taxa above contain freshwater and/or terrestrial groups. These are expected to occur in such sedimentary environments, which is recorded in the `environment` field. We can list the sampled environments with a single line of code.

```
# filter by environment
levels(factor((dat$environment)))
```

The following environments are expected to contain terrestrial taxa.

```
omitEnv <- c(
  "\"floodplain\"", "alluvial fan", "cave", "\"channel\"", "channel lag" ,
  "coarse channel fill", "crater lake", "crevasse splay", "dry floodplain",
  "delta plain", "dune", "eolian indet.", "fine channel fill", "fissure fill",
  "fluvial indet.", "fluvial-lacustrine indet.", "fluvial-deltaic indet.",
  "glacial", "interdune", "karst indet.", "lacustrine - large",
  "lacustrine - small", "lacustrine delta front", "lacustrine delta plain",
  "lacustrine deltaic indet.", "lacustrine indet.",
  "lacustrine interdistributary bay", "lacustrine prodelta", "levee", "loess",
  "mire/swamp", "pond", "sinkhole", "spring", "tar", "terrestrial indet.",
  "wet floodplain")
```

We can omit occurrences with these entries with a similar command as above.

```
# omit the occurrences
dat <- dat[!dat$environment%in%omitEnv, ]
```

Although this filtering is not perfect, we believe that it is adequate for answering large-scale questions. The remaining non-marine and non-animal occurrences are unlikely to influence the results. After this filtering step

```
nrow(dat)
```

```
## [1] 880576
```

occurrences remain in the dataset. Collections that came from unlithified sediments can yield fossils with unusually good preservation. As these are more frequent in younger sites and occur heterogeneously, sampling bias can be reduced by omitting such collections from the data.

```
dat <- dat[dat$lithification1!="unlithified", ]
```

This last step leaves

```
nrow(dat)
```

```
## [1] 815383
```

occurrences for the analyses.

3.3. Stratigraphic binning

To use the occurrences in the `divDyn` package in an efficient way, the occurrence entries have to be assigned to a discrete time scale. We included two of these time-scales in the package: the widely-used 10 myr time scale of the Paleobiology Database, and another one based on the stratigraphic stages of Ogg et al. (2016). These can be loaded with the `data()` function.

```
# 10 million year timescale
data(bins)
# stage-level timescale
data(stages)
```

You can learn more about these time scales if you type in `?bins` or `?stages`. The first time scale (`bins`) has 49 entries identifying roughly 10-million year bins. The second one (`stages`) has almost double the stratigraphic resolution, but some of the ICS stages are clumped to ensure a more even distribution of durations (cf. Miocene and Pliocene).

It is easier to handle the two timescales with the same functions if the time bin names have identical column names in both tables.

```
# names of the bins
colnames(bins)[colnames(bins)=="X10"] <- "name"
# names of the bins
colnames(stages)[colnames(stages)=="stage"] <- "name"
```

The original data in the Paleobiology Database get their stratigraphic information based on the `early_interval` and `late_interval` values. The valid entries in these variables come from a list of interval names that convert them to numeric ages and establish connections between the different entries. In the dynamic timescale of Fossilworks (J. Alroy), these entries are also linked to the 10 million year timescale and the ICS stages (Ogg et al. 2016), without the changes in the Neogene. The `early_interval` and `late_interval` values designate stratigraphic position in a straightforward way: the `early_interval` marks the oldest possible age and the `late_interval` the youngest. If a single name suffices to describe the inherent uncertainty of an interval, the `late_interval` remains empty.

Using a complete download of collections from Fossilworks, we compiled a table that resolves these ‘interval’ entries to the timescales of our interest, which can be viewed in the `stratkeys` object, which can be invoked by `data(stratkeys)`.

The entries in this table were then transformed to `list` type entries to enable a grouping function `categorize()` that replaces groups of entries in a vector with single group names (see `?categorize`). The `keys` object contains the information in a relevant form and includes additional variables to group occurrences based on environmental and lithological information.

```
data(keys)
# using the others will be included in an appendix to this vignette
names(keys)
```

```
## [1] "binInt" "stgInt" "reefs" "lith" "lat" "bath" "grain"
```

The stratigraphic binning starts with figuring out which numbered bins the `early_interval` and `late_interval` entries are assigned to. Let’s start with the 10 million year bins.

```
# categorize entries as they are in the lookup table
binMin <- categorize(dat[, "early_interval"], keys$binInt)
binMax <- categorize(dat[, "late_interval"], keys$binInt)
```

Then the entries have to be converted to simple numeric values.

```
binMin <- as.numeric(binMin)
binMax <- as.numeric(binMax)
```

This code creates two vectors of numeric bin numbers, where NA entries indicate that the names were not found in the table, and -1 entries indicate empty character strings – where no `late_interval` entry is given. As our goal is to retain only those occurrences that have precise enough stratigraphic assignments, we only want to consider those occurrences that have either the same `binMin` or `binMax` number or where `binMax` is -1. This is accomplished by the following steps.

First, a final, empty vector is defined.

```
dat$bin <- rep(NA, nrow(dat))
```

Then the condition above is expressed indicating which rows have only a single assigned bin number.

```
binCondition <- c(
  # the early and late interval fields indicate the same bin
  which(binMax==binMin),
  # or the late_interval field is empty
  which(binMax==-1))
```

Finally, those values are copied, where the condition is true.

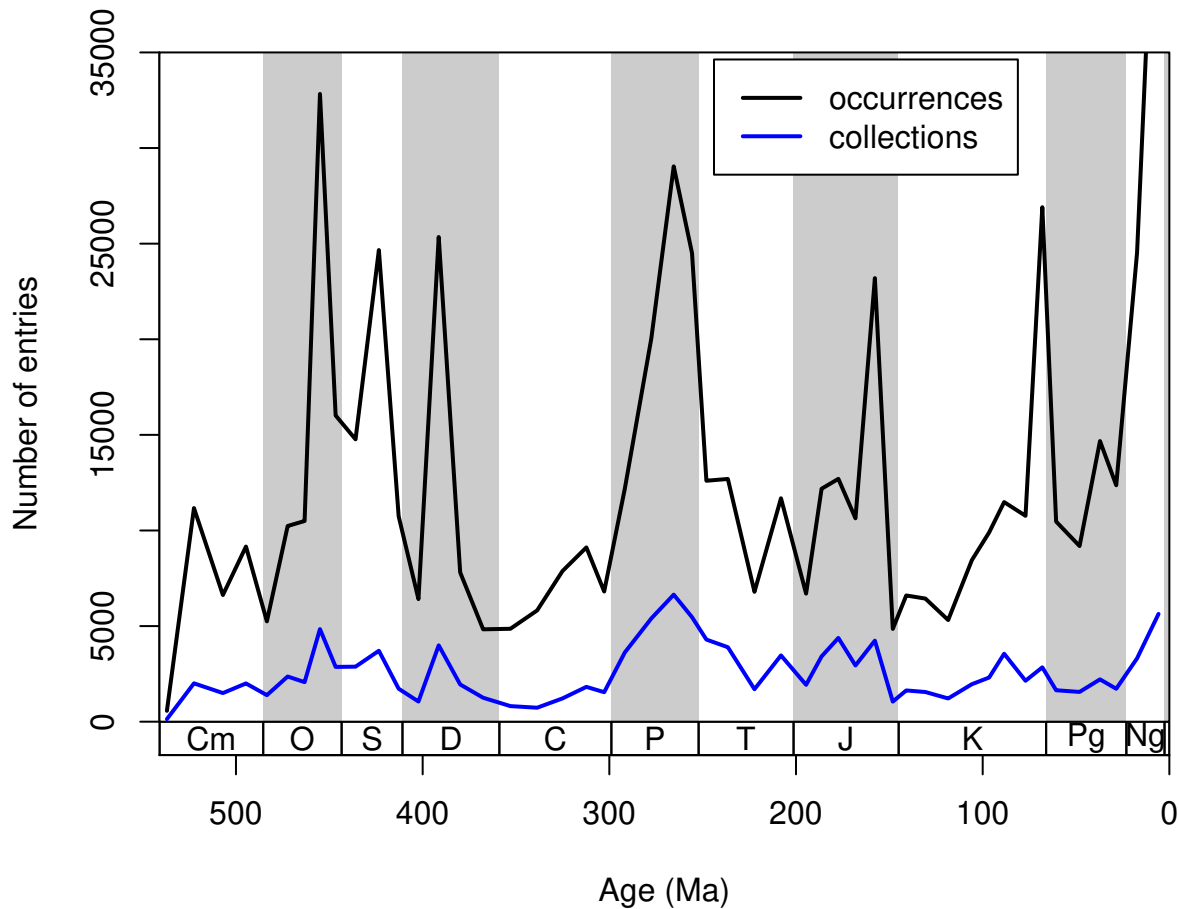
```
# in these entries, use the bin indicated by the early_interval
dat$bin[binCondition] <- binMin[binCondition]
```

The final object is a column in the data table, where `dat$bin` is a single variable of integers that have NA entries where the collection/occurrence cannot be assigned to a single bin in the time scale. The number of occurrences and collections in each bin can be calculated with `table()` or in a single step with the `binstat()` function.

```
sampBin <- binstat(dat, tax="clgen", bin="bin", coll="collection_no", duplicates=FALSE)
```

As every row in `dat` represents a species-level occurrence entry, multiple species of the same genus can be registered in every collection. These are only counted as one occurrence, when the `duplicates` argument of `binstat()` is set to `FALSE`. The resulting `data.frame` can be used to plot the number of occurrences and collections through time.

```
# the plot
tsplot(stages, boxes="per", shading="per", xlim=4:95, ylim=c(0,35000),
  ylab="Number of entries", xlab="Age (Ma)")
# occurrences
lines(bins$mid, sampBin$occs, lwd=2)
# collections
lines(bins$mid, sampBin$colls, lwd=2, col="blue")
# legend
legend("topright", bg="white", legend=c("occurrences", "collections"),
  col=c("black", "blue"), lwd=2, inset=c(0.15,0.01), cex=1)
```

The following code repeats the entire process for the stratigraphic stages:

```
# the 'stg' entries (lookup)
stgMin <- categorize(dat[, "early_interval"], keys$stgInt)
stgMax <- categorize(dat[, "late_interval"], keys$stgInt)

# convert to numeric
stgMin <- as.numeric(stgMin)
stgMax <- as.numeric(stgMax)

# empty container
dat$stg <- rep(NA, nrow(dat))

# select entries, where
stgCondition <- c(
  # the early and late interval fields indicate the same stg
  which(stgMax==stgMin),
  # or the late_interval field is empty
  which(stgMax==1))

# in these entries, use the stg indicated by the early_interval
```

```
dat$stg[stgCondition] <- stgMin[stgCondition]
```

But beware! Stage-level assignments have been a problem in the earliest Paleozoic (Cambrian and Ordovician periods). Therefore, the assignments above are not perfect in these periods: numerous entries are not properly processed based on the interval lookup table. To make use of these collections in the analyses that follow, we processed these data separately. The necessary files to assign these corrections are added to the example GitHub repository and will be updated for future use. Cambrian collections were assigned one-by-one by Na Lin for the analysis of studying the diversity dynamics of the Cambrian (Na and Kiessling, 2015). You can download these from:

```
load(url(
  "https://github.com/adamkocsis/ddPhanero/raw/master/data/Stratigraphy/2018-08-31/cambStrat.RData"))
```

And apply them with the following script:

```
source(
  "https://github.com/adamkocsis/ddPhanero/raw/master/scripts/strat/2018-08-31/cambProcess.R")
```

For the Ordovician assignments, Wolfgang Kiessling compiled tables on formations and biozones.

```
load(url(
  "https://github.com/adamkocsis/ddPhanero/raw/master/data/Stratigraphy/2018-08-31/ordStrat.RData"))
```

You can use these with the following script.

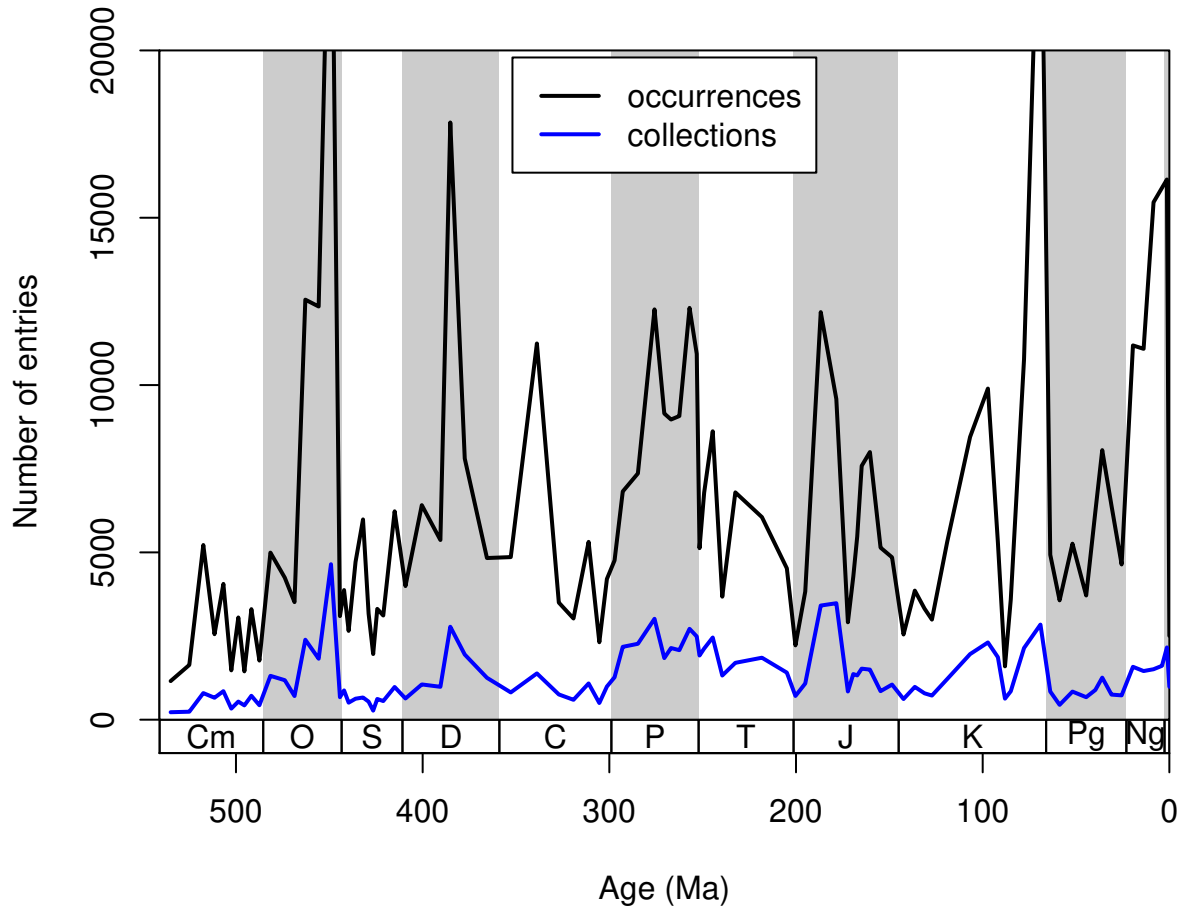
```
source(
  "https://github.com/adamkocsis/ddPhanero/raw/master/scripts/strat/2018-08-31/ordProcess.R")
```

Now that we the occurrences assigned to stages, as above, the number of sampled occurrences and collections can be calculated with `binstat()`.

```
sampStg <- binstat(dat, tax="clgen", bin="stg", coll="collection_no", duplicates=FALSE)
```

And then can be plotted in a similar way.

```
# the plot
tsplot(stages, boxes="per", shading="per", xlim=4:95, ylim=c(0,20000),
  ylab="Number of entries", xlab="Age (Ma)")
# occurrences
lines(stages$mid, sampStg$occs, lwd=2)
# collections
lines(stages$mid, sampStg$colls, lwd=2, col="blue")
# legend
legend("top", bg="white", legend=c("occurrences", "collections"),
  col=c("black", "blue"), lwd=2, inset=c(0.15,0.01), cex=1)
```



4. Calculating richness (Figure 1)

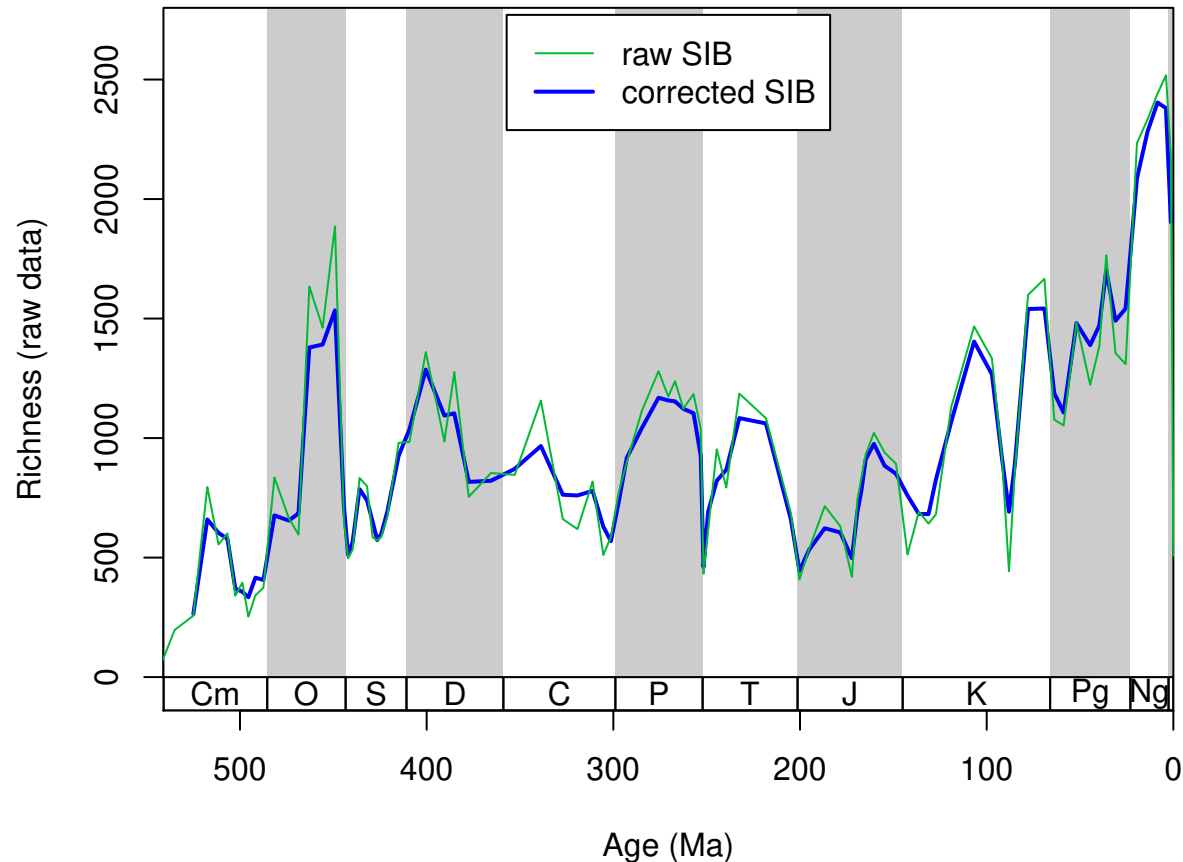
The first figure of the paper depicts changes in genus richness over the Phanerozoic, using the 10 million year intervals with raw and subsampled data. You can calculate the raw results with the basic `divDyn()` function.

```
# raw patterns
ddStages <- divDyn(dat, bin="stg", tax="clgen")
```

This function call produces a `data.frame` class object, with the calculated metrics as columns. The variables are explained in the function help file `?divDyn`. Because we have been using positive integer bin numbers, the indices of rows in this table match the bin identifier numbers. Therefore, the corresponding age values are in the same index rows in the `bins` table. To visualize the results, we need to plot the timescale first with the `tsplot()` function (otherwise it is difficult to match results with time intervals visually) and then we can use `lines()` to show the variable in question.

```
tsplot(stages, boxes="per", shading="per", xlim=4:95, ylim=c(0,2800),
       ylab="Richness (raw data)", xlab="Age (Ma)")
lines(stages$mid, ddStages$divCSIB, col="blue", lwd=2)
lines(stages$mid, ddStages$divSIB, col="#00BB33", lwd=1)
```

```
legend("top", inset=c(0.01,0.01), legend=c("raw SIB", "corrected SIB"),
      col=c("#00BB33", "blue"), lwd=c(1,2), bg="white", cex=1)
```



Sampling-standardized values can be calculated with a single function called `subsample()`.

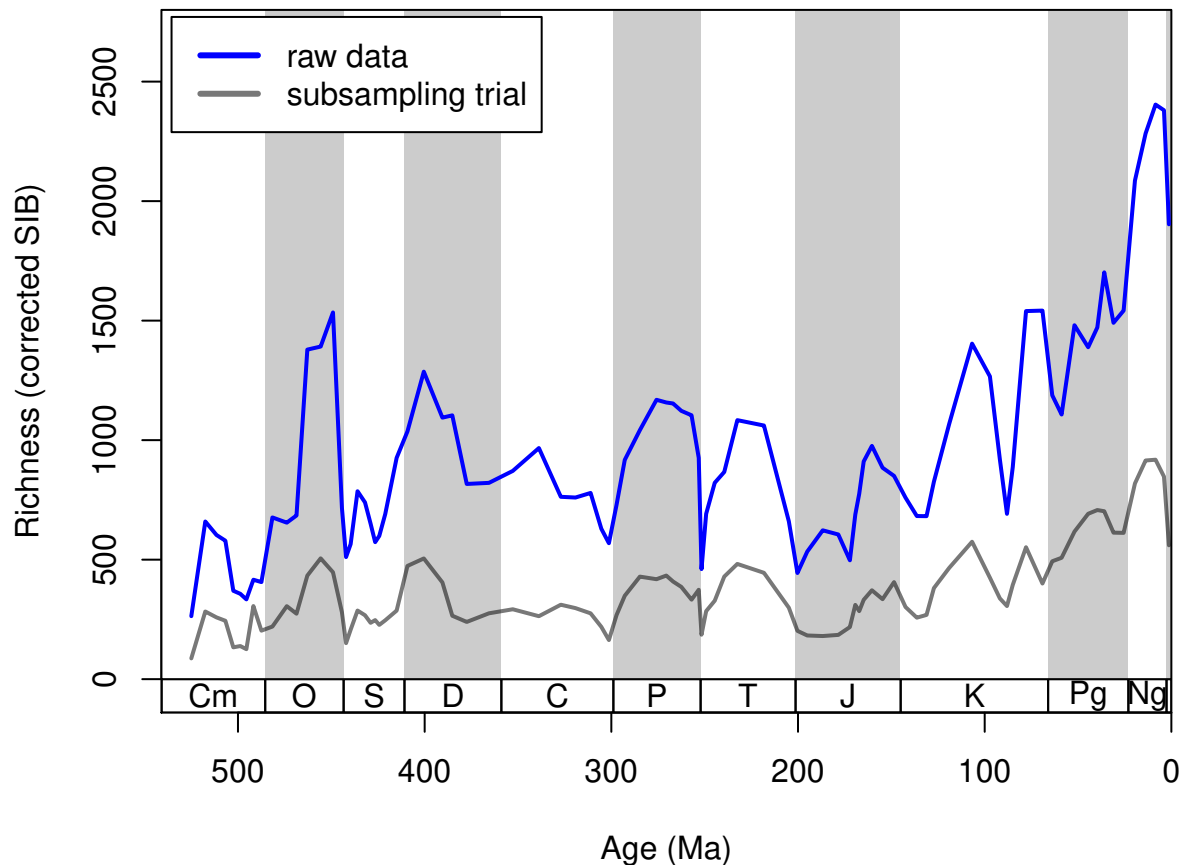
```
sqsStagesPlot <- subsample(dat, bin="stg", tax="clgen", coll="collection_no", q=0.7,
  iter=100, ref="reference_no", singleton="ref", type="sqs", duplicates=FALSE,
  excludeDominant=TRUE, largestColl =TRUE, output="dist")
```

The function above is configured to use the ‘Shareholder Quorum Subsampling’ (Alroy, 2010a, 2010b) or ‘coverage-based rarefaction’ (Chao and Jost, 2012) method (`type="sqs"`) that subsamples the data down to an even level of sample coverage (Good, 1953). The former name is usually used to refer to the algorithmic version of the method; latter one to refer to the analytical version. We only implemented the algorithmic version as the complexities of paleontological data make the analytical version less useful in this context. The function is configured to use a reference-based ‘singleton’ treatment (`singleton="ref"`) for overall sampling correction, excluding dominant taxa from all calculations involving frequencies (`excludeDominant=TRUE`) and with the separate treatment of the largest collection in the time slice (`largestColl=TRUE`), as indicated by Alroy (2010a). Setting `output` to `"dist"` makes the function return the results of individual subsampling trials. Please take a look at the help files of the functions `?subsample` or `?subtrialsSqs` if you want to know more.

You can compare one trial result (in this case the 51st) with the original time series with the following code:

```
tsplot(stages, boxes="per", shading="per", xlim=4:95, ylim=c(0,2800),
      ylab="Richness (corrected SIB)", xlab="Age (Ma)")

lines(stages$mid, ddStages$divCSIB, col="blue", lwd=2)
lines(stages$mid, sqsStagesPlot$divCSIB[,51], col="#00000088", lwd=2)
legend("topleft", bg="white", legend=c("raw data", "subsampling trial"),
      col=c("blue", "#00000088"), lwd=3, inset=c(0.01, 0.01))
```

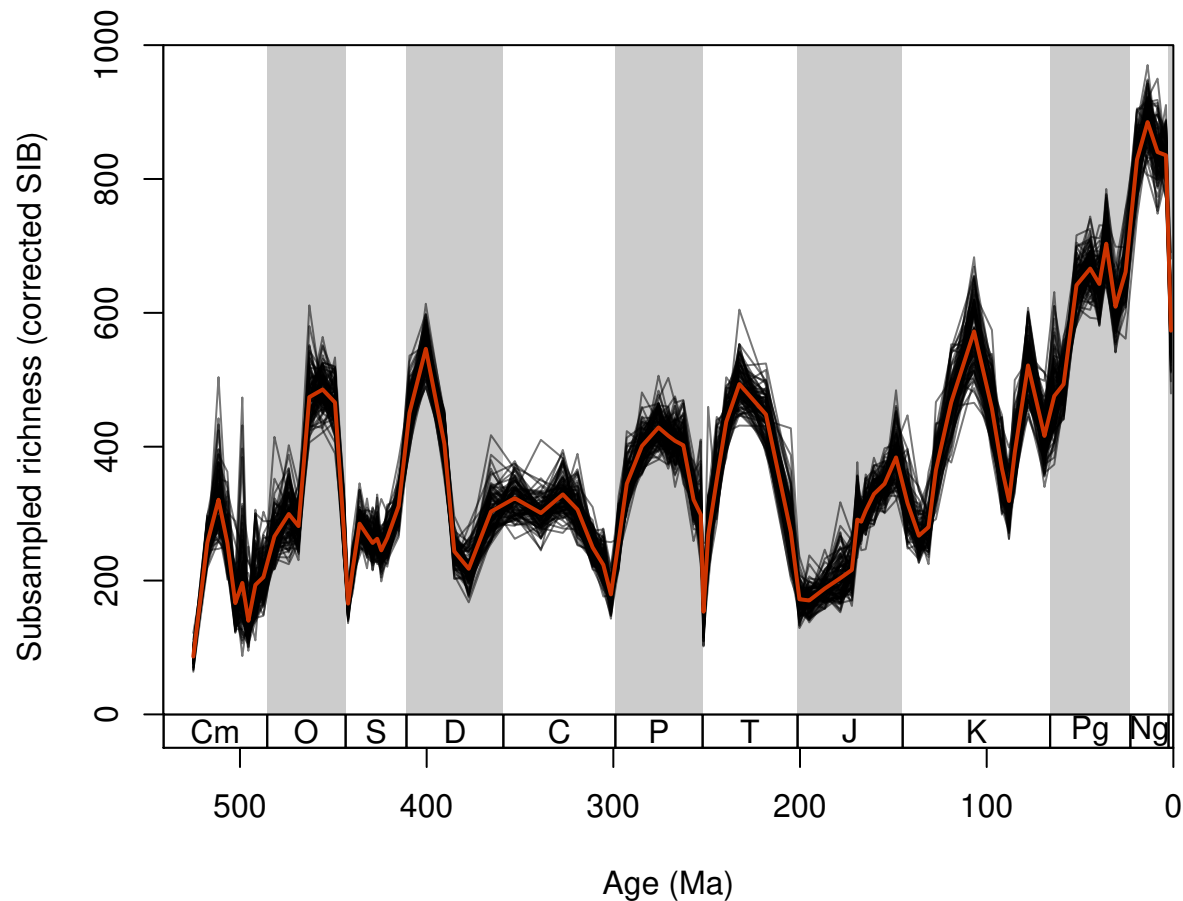


You can either let the function average the trial results by setting the `output` argument accordingly, or you can do it yourself. Here are the results of the 100 trials calculated above.

```
# only plot
tsplot(stages, boxes="per", shading="per", xlim=4:95, ylim=c(0,1000),
      ylab="Subsampled richness (corrected SIB)", xlab="Age (Ma)")

# loop through all trial results
for(i in 1:ncol(sqsStagesPlot$divCSIB)){
  lines(stages$mid, sqsStagesPlot$divCSIB[,i], col="#00000088")
}
```

```
# the mean of the trial results
meanRes <- apply(sqsStagesPlot$divCSIB, 1, mean, na.rm=T)
lines(stages$mid, meanRes, col="#CC3300", lwd=2)
```



5. Analyzing the taxonomic rates

The analyses in this study focus on turnover rates rather than richness estimates. The prototypes for these are the analyses presented by Alroy (2008). We provide the analytical code to reproduce these using the per capita rates (Foote, 1999). These are the most used in the literature, even though they were critiqued by Alroy (2014) for their accuracy. The stratigraphic resolution in this is example was at the level of stages. All necessary functions are included in the phanDyn.R file, which is deposited on GitHub.

```
# load the necessary functions
source("https://github.com/adamkocsis/ddPhanero/raw/master/scripts/0.3/phanDyn.R")
```

5.1. Questions addressed with original rate values

The analyses in the following part of this section can be repeated with any set of extinction rate, origination rate and richness time series. To prepare the generalized implementation of the following code in a function (Section 5.3), the necessary variables are renamed so the same analyses can be rerun on different series if they are named appropriately. Variables from the timescale tables are also necessary.

```
# the name of the overall data frame (already created)
ext <- ddStages[ , "extPC"] # extinction rates
ori <- ddStages[ , "oriPC"] # origination rates
div <- ddStages[ , "divCSIB"] # diversity (richness) series
dur <- stages[ , "dur"] # durations of time intervals
age <- stages[ , "mid"] # time interval midpoints
name <- stages[ , "name"] # names of the time intervals
```

A. Is the pulsed model of turnover supported by the instantaneous rates?

The taxonomic rates currently implemented in the package converge on the per capita rates (Foote, 1999) as sampling completeness approaches 1. The original definition of the per capita rates expresses the intensity of turnover as per lineage-myr, which means that the log proportions are divided (normalized) by the durations of the time intervals. Implicit in these equations is that taxonomic turnover happened continuously in the bin, which has been suggested to be an invalid assumption, at least for extinctions (known as the pulsed model: Foote, 2005). If a ‘pulsed’ model is correct, the normalization of rates by bin duration is not only unnecessary but plain wrong. A ‘pulsed’ model is supported by the absence of correlations between turnover rates and bin durations (Alroy, 2008). With normalization a negative correlation between rate and duration may arise (Alroy, 2008). Rates calculated this way are sometimes characterized as ‘instantaneous’.

Using the `pulseCont()` function we prepared tests for correlations between interval durations and rate values with and without the normalization. It uses tests of Spearman’s rank-order correlations to assess these relationships. The function can be found in the `phanDyn.R` file, and it takes a rate series and vector of bin durations as arguments.

```
pulseCont(ext, dur, alpha=0.01)

## $est
## not-normalized      normalized
##      0.2311910      -0.4545859
##
## $p
## not-normalized      normalized
##  2.835024e-02      6.767625e-06
##
## $sig
## normalized
## -0.4545859
```

The `est` element of the output list contains the correlation coefficients, `p` the p -values, and the `sig` element indicates significant results at the $p < 0.01$ level. As only the normalized rates are significantly correlated with bin durations, it is likely that normalization strengthens the association between the two variables. This suggests that for extinctions (based on the raw, stage-level per capita rates) the ‘pulsed’ model turnover is more appropriate.

B. Are the taxonomic rates declining?

The decline of taxonomic turnover rates is one of the most robust findings of paleobiology and has been reported many times in the literature (e.g. Raup and Sepkoski, 1982; Sepkoski, 1998, Bambach et al, 2004; Alroy, 2008). Therefore, we expect both origination and extinction rates to decline over time. In general, we have to be cautious about the earliest and latest parts of the time series. First, these rates cannot be estimated for the first and last bins (some even not for the second and the second last), as it is impossible to count taxa with certain temporal patterns of occurrences in these intervals (e.g. three-timers). The so-called edge-effects (Foote 2000) also bias the per-capita rates systematically. The stratigraphy of the Cambrian and the Ordovician intervals remains problematic, while we also know that higher taxonomic entities originated in these intervals that separate them from the rest of the Phanerozoic. Therefore, we decided to copy the analytical parameters of Alroy (2008), and the declining pattern was also assessed separately for the post-Ordovician interval. The following intervals were applied in the assessment:

```
# total decline interval in Ma
maxAgeDecline <- 540
minAgeDecline <- 20
# mid dates of post-Ordovician time bins (both time scales) are younger than
postDate <- 443
```

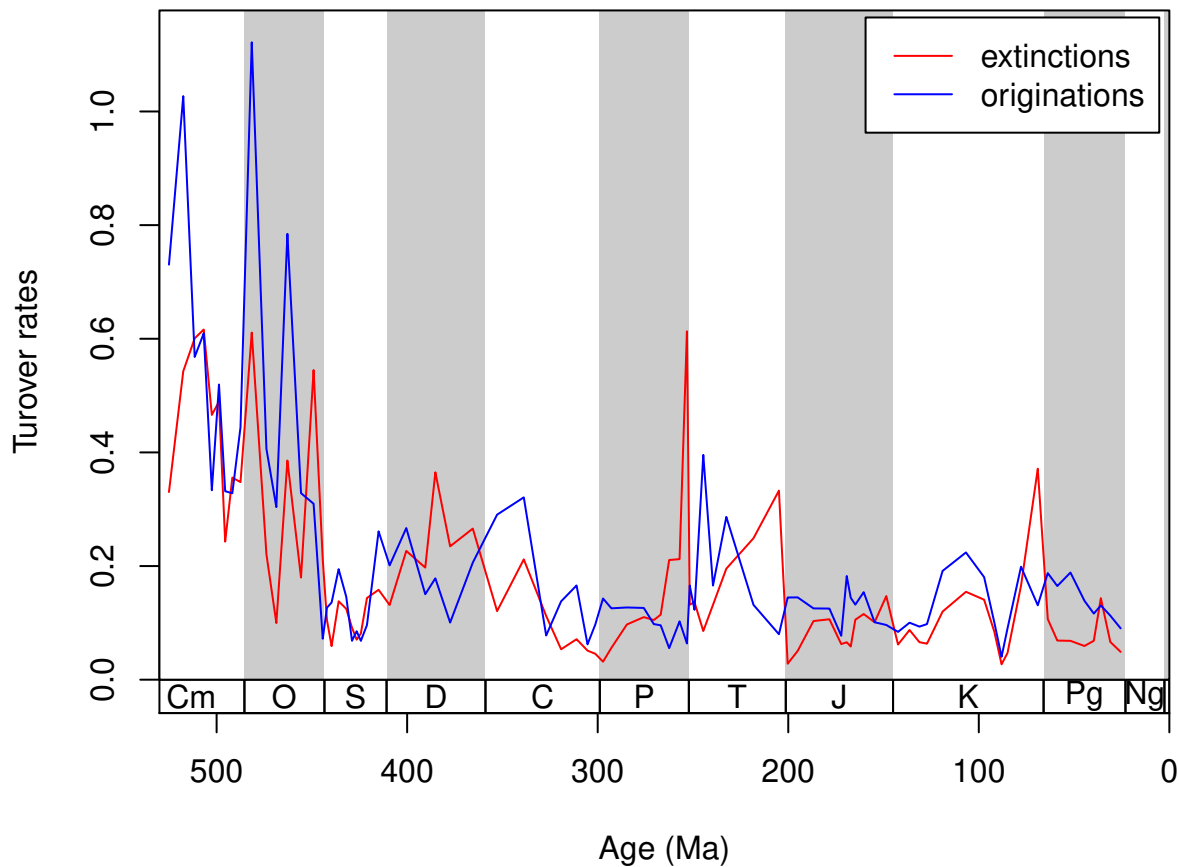
To retain the match between the bin numbers in the time scale table and the indices of the results, we defined logical vectors that replace values that are out of this interval with NA entries.

```
# logical vectors indicating what is necessary
# Cambrian included
indNAAlong <- !(age<=maxAgeDecline & age>=minAgeDecline)
# post Ordovician
indNApostord <- !(age<=postDate & age>= minAgeDecline)
```

You can relate these values to the estimated ages in **bins** and **stages** objects. It is good practice to visualize the variables before carrying out the correlation tests.

```
# actual variables used in the correlation tests:
extVarLong <- ext
oriVarLong <- ori
extVarLong[indNAAlong] <- NA
oriVarLong[indNAAlong] <- NA

# plotting
# maximum y value
nMaxRate<-max(c(extVarLong, oriVarLong), na.rm=T)
# actual plot
tsplot(stages, boxes="per", shading="per", ylim=c(0, nMaxRate*1.05),
       ylab="Turnover rates", xlim=c(530,0), xlab="Age (Ma)")
lines(age, extVarLong, col="red")
lines(age, oriVarLong, col="blue")
legend("topright", legend=c("extinctions", "originations"), col=c("red", "blue"),
       lwd=c(1,1), bg="white", inset=c(0.01, 0.01))
```

Some patterns of a decline are evident from viewing these series, but we need to test them. We can create a 2 by 4 matrix to hold results of correlations against time (coefficients and *p*-values) that will be calculated for each of the series.

```
# create a table from the decline information
decMat <- matrix(NA, ncol=2, nrow=4)
colnames(decMat)<- c("correlation","p-value")
rownames(decMat)<- c("extinction", "post-Ordovician extinction",
"origination", "post-Ordovician origination")
```

First, the total declines (extVarLong and oriVarLong) are assessed.

```
# is the extinction rate declining
extDecline <- cor.test(age, extVarLong, method="spearman")
decMat[1, 1] <- extDecline$estimate
decMat[1, 2] <- extDecline$p.value

# is the origination rate declining
oriDecline <- cor.test(age, oriVarLong, method="spearman")
decMat[3, 1] <- oriDecline$estimate
decMat[3, 2] <- oriDecline$p.value
```

Then the post-Ordovician declines are assessed.

```
# extinction
extPostOrd <- extVarLong
extPostOrd[indNApostord] <- NA
extDeclinePostOrd <- cor.test(age, extPostOrd, method="spearman")
# save
decMat[2, 1] <- extDeclinePostOrd$estimate
decMat[2, 2] <- extDeclinePostOrd$p.value

# origination
oriPostOrd <- oriVarLong
oriPostOrd[indNApostord] <- NA
oriDeclinePostOrd <- cor.test(age, oriPostOrd, method="spearman")
# save
decMat[4, 1] <- oriDeclinePostOrd$estimate
decMat[4, 2] <- oriDeclinePostOrd$p.value
```

The results of the test above are in the `decMat` object, clearly suggesting patterns of an overall decline for both extinctions and originations, although this might not be true for the post-Ordovician parts of the series.

```
round(decMat, 4)
```

##	correlation	p-value
## extinction	0.5213	0.0000
## post-Ordovician extinction	0.2309	0.0564
## origination	0.4449	0.0000
## post-Ordovician origination	0.0602	0.6221

5.2. Detrending the series

The definition of the rate distribution, identification of outlying values and implementation of accurate tests of cross correlations between the time series assume that the series are stationary. To get closer to meeting this assumption, the rate series have to be detrended. There are multiple ways of fitting a trend to the time series. You can fit a polynomial function, model it as an ARIMA process, or apply a scatterplot smoother (e.g. LOESS) to the values.

Among these options, the latter is the simplest and uses the lowest number of assumptions. This is a very useful property, as it decreases the chance of inappropriate fit when the same analytical code is run on different time series, or when long-term trends are not significant. This section starts by going through how this procedure can be applied to the time series of originations and extinctions without transforming the data (see section 5.3 for other options).

LOESS relies on a bandwidth parameter (or `span`) that describes how much smoothing should be applied to the series, which is usually considered to be an arbitrary parameter. Luckily, there are ways to search for optimal setting of this argument. An easily applicable implementation is included in the `fANCOVA` package (Wang, 2010), which you can download from the CRAN servers with the following line of code:

```
install.packages("fANCOVA")
```

The relevant function requires the time series to be free of missing values with the corresponding ages. You can get rid of these in a structured way by executing these commands:

```
# extinctions
extMiss <- !is.na(extVarLong)
transExtNoNA <- extVarLong[extMiss]
ageExt <- age[extMiss]
```

```
# originations
oriMiss <- !is.na(oriVarLong)
transOriNoNA <- oriVarLong[oriMiss]
ageOri <- age[oriMiss]
```

After running this snippet, `transExtNoNA` includes the extinction rate values and `ageExt` has the corresponding ages. The same applies to `transOriNoNA` and `ageOri` but for originations. Doing the actual smoothing is just a single step with `fANCOVA`.

```
# the models
extModel <- fANCOVA::loess.as(ageExt, transExtNoNA, degree = 1,
  criterion = "aicc", user.span = NULL, plot = FALSE)
oriModel <- fANCOVA::loess.as(ageOri, transOriNoNA, degree = 1,
  criterion = "aicc", user.span = NULL, plot = FALSE)
```

These settings will make the function use corrected AICs to define the optimal smoothing of the curves. To get the actual values, predictions must be calculated from these model objects. Running the `predict()` function on them will output the predicted values (i.e. the trend).

```
# predictions for extinctions
predict(extModel)
```

```
## [1] 0.53746494 0.50668275 0.48211034 0.46260995 0.44507707 0.42950691
## [7] 0.41611768 0.40095857 0.38376119 0.36006301 0.32801215 0.30557572
## [13] 0.28073744 0.25133239 0.22637535 0.21011263 0.20440120 0.20006137
## [19] 0.19372742 0.18935600 0.18866048 0.18726434 0.18537566 0.18353155
## [25] 0.18244418 0.18350675 0.18681380 0.18842407 0.18749597 0.18403660
## [31] 0.17419629 0.15814605 0.14158224 0.13354997 0.13103013 0.13072435
## [37] 0.13187529 0.12974943 0.12516314 0.12500616 0.13147764 0.13745002
## [43] 0.14355639 0.14839663 0.15333175 0.15874459 0.16277745 0.16448201
## [49] 0.16711704 0.17040721 0.17183631 0.17044740 0.15817089 0.14297113
## [55] 0.13807548 0.13262001 0.12414259 0.11558811 0.10943667 0.10691239
## [61] 0.10530732 0.10373145 0.10144512 0.09931574 0.09793809 0.09764963
## [67] 0.09834304 0.09964924 0.10104551 0.10366881 0.10618414 0.10807744
## [73] 0.10761853 0.10648701 0.10610469 0.10587751 0.10458709 0.10356794
## [79] 0.10260090 0.10152715 0.10019635 0.09912042 0.09842842 0.09723028
## [85] 0.09560299
```

There is a problem with this result, however, namely that certain predictions where the original time series had NA values are not present. This leads to misalignment with the other time series and numerous future problems. Therefore, it is better practice to coerce the output of missing values in the prediction, by specifying the `x` (time) coordinates for the predictions. These are in the original `age` vector and can be used this way:

```
predExt <- predict(extModel, newdata=data.frame(x=age)) # extinctions
predOri <- predict(oriModel, newdata=data.frame(x=age)) # originations
```

The resulting objects have the same number of values as the original series.

```
length(predExt)
```

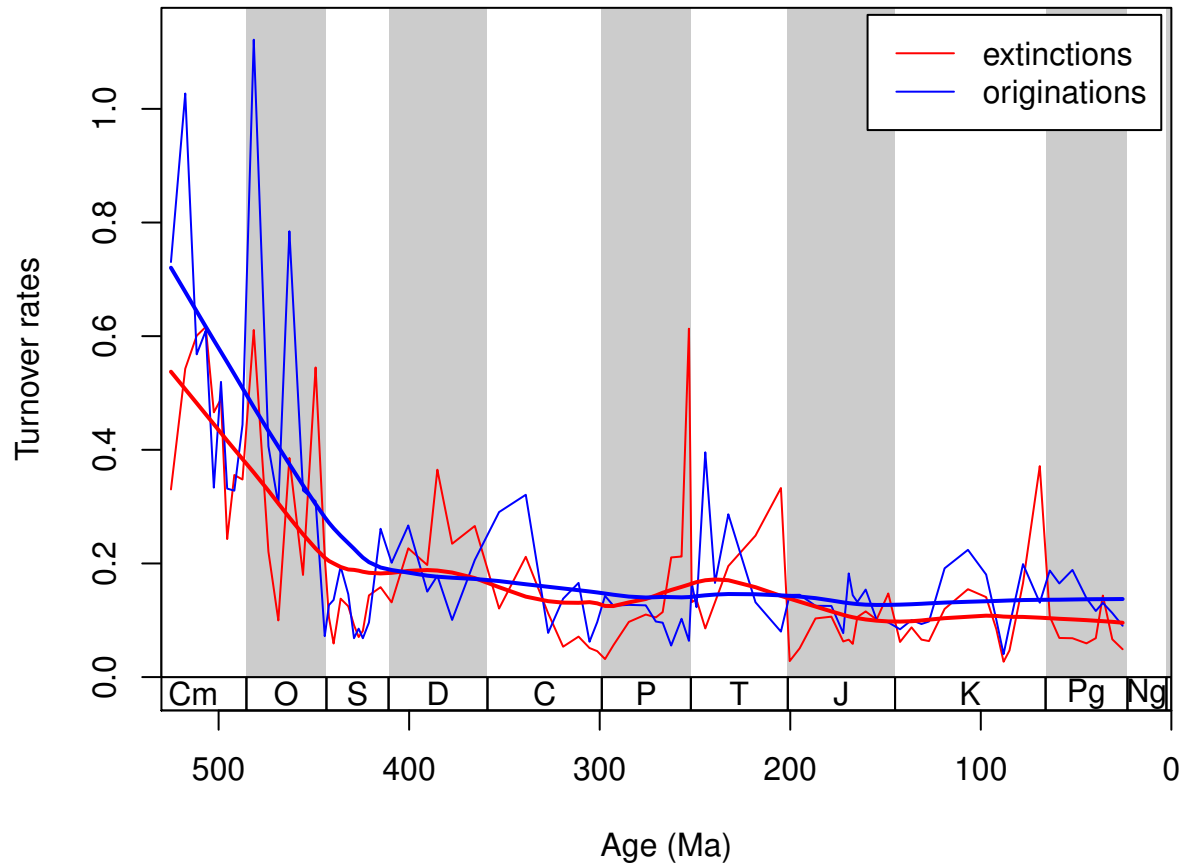
```
## [1] 95
```

```
length(predOri)
```

```
## [1] 95
```

You can draw the predicted values with `lines()`.

```
lines(age, predExt,col="red", lwd=2) # extinctions
lines(age, predOri,col="blue", lwd=2) # originations
```



The detrended values are then the residuals with multiplicative decomposition.

```
detExt <- extVarLong/predExt
detOri <- oriVarLong/predOri
```

Following Alroy (2008), only the post-Cambrian values were used in further analyses, as the Cambrian rates appear to be unusually high in all cases.

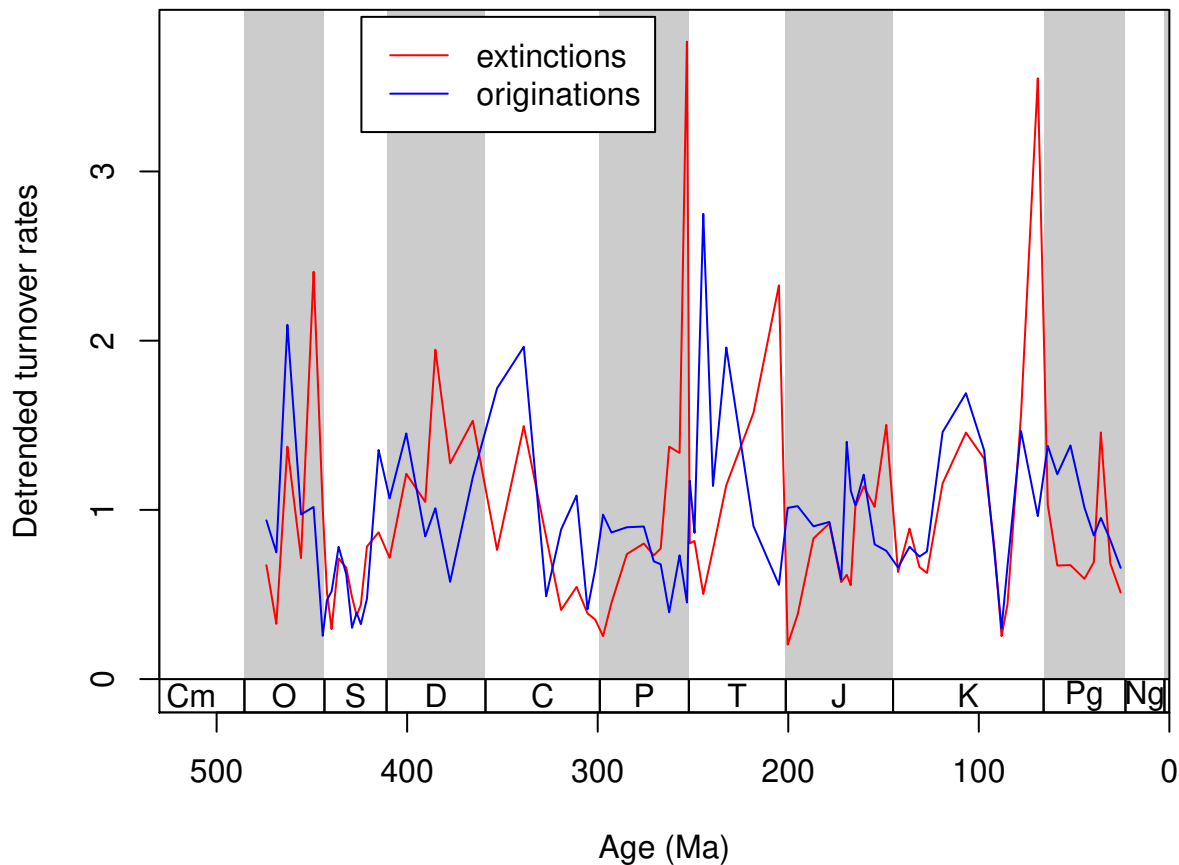
```
# interval for detrending in ma
maxAgeDetrend <- 475
minAgeDetrend <- 20
# indices and variables
indNashort<- !(age<=maxAgeDetrend & age>=minAgeDetrend)

# recreate the variables
detExtShort <- detExt
detOriShort <- detOri
detExtShort[indNashort] <- NA
```

```
detOriShort[indNAsShort] <- NA
```

You can plot the residuals with `lines()`.

```
# for setting the y-axis range
detMax <- max(c(detExtShort, detOriShort), na.rm=T)
# plot
tsplot(stages, boxes="per", shading="per", ylim=c(0, detMax*1.05),
       ylab="Detrended turnover rates", xlim=c(530, 0), xlab="Age (Ma)")
lines(age, detExtShort, col="red") # extinctions
lines(age, detOriShort, col="blue") # originations
legend("topleft", legend=c("extinctions", "originations"), col=c("red", "blue"),
       lwd=c(1, 1), bg="white", inset=c(0.2, 0.01))
```

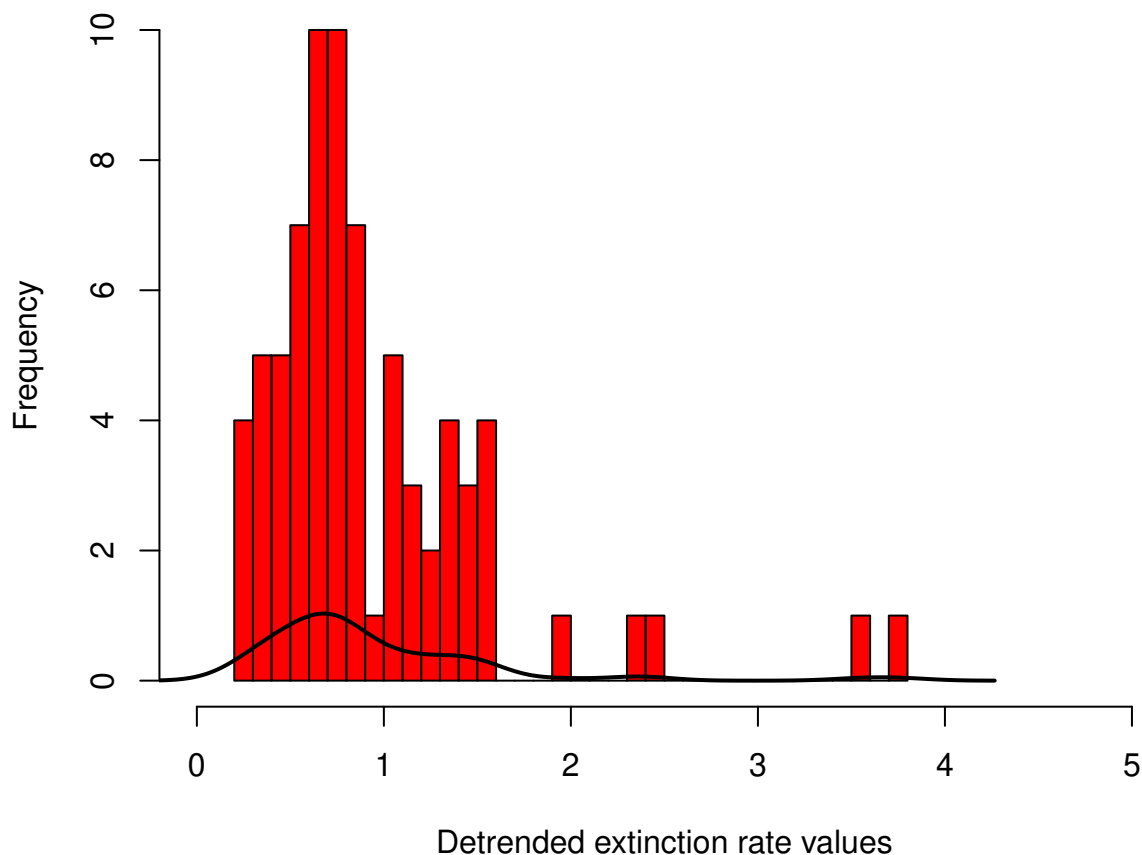


As mentioned before, there are other ways to do the detrending of the series, some of which we also implemented (see section 5.3 for options). You are invited to inspect these within the `analyzedMetrics()` function that can be found in the `phanDyn.R` file.

C. Is the distribution of rates lognormal?

Mass extinctions have always been diagnosed by analyzing the distributions of the rates. Originally, these were just compared to the point estimation confidence-interval of a linear decline model (Raup and Sepkoski, 1982), which was already criticized back then (Quinn, 1983). Bambach et al. (2004) fitted a LOESS model to the extinction rates, and analyzed the distribution of its residuals (similar to what we just did, assuming a local background process). Alroy (2008) suggested detrending the rates with exponential functions. He analyzed the distribution of the rates first and then pointed to potential outliers. Here is the distribution of the detrended extinction rate series with a kernel density estimator:

```
# extinctions histogram
hist(detExtShort , breaks=30, xlim=c(0,5), col="red",
     xlab="Detrended extinction rate values", main="")
# kernel density estimator
den <- density(detExtShort [!is.na(detExtShort )])
lines(den$x, den$y, lwd=2)
```



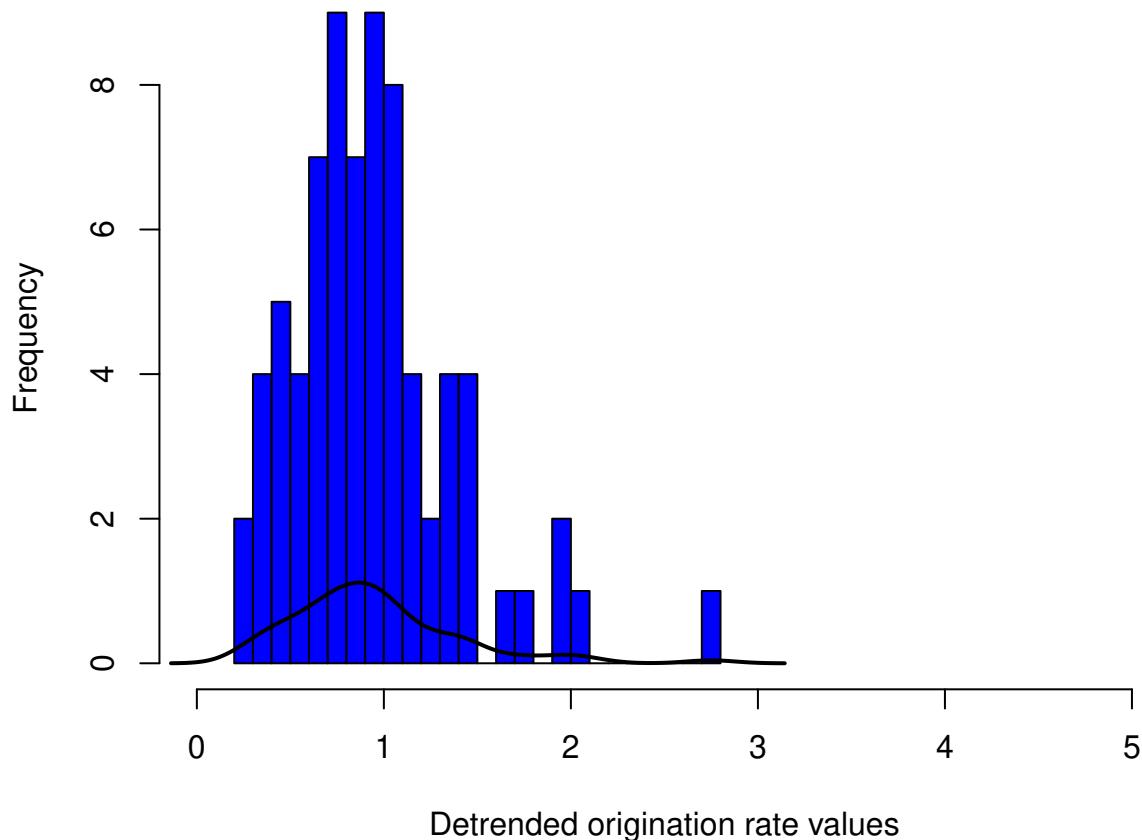
And here is the same thing for originations.

```
# histogram
hist(detOriShort, breaks=30, xlim=c(0, 5), col="blue",
```

```

xlab="Detrended origination rate values", main="")
# kernel density estimator
den<-density(detOriShort[!is.na(detOriShort)])
lines(den$x, den$y, lwd=2)

```



As the extinction rate distributions tend to be right-skewed and they are 0-bounded, analyzing the logarithms of the rates, or their square roots make more sense. The question behind the analysis of rate distributions is whether we can distinguish between two processes of extinction in terms of mass extinction episodes and background extinctions, or whether the two just at different positions along a spectrum. A lognormal distribution suggests that there is only a difference in magnitude between these intervals, which has to be assessed. Therefore, the function checks whether the rates can come from a lognormal distribution or not, which is implemented by using Shapiro-Wilk tests.

```

# are logged data normal? - extinction
extVar <- log(detExtShort)
extVar[is.infinite(extVar)] <- NA # omit infinities
pShap <- shapiro.test(extVar)$p.value
names(pShap)[1] <- "ext"

# are logged data normal? - origination

```

```
oriVar <- log(detOriShort)
oriVar[is.infinite(oriVar)] <- NA # omit infinities
pShap <- c(pShap, shapiro.test(oriVar)$p.value)
names(pShap)[2] <- "ori"

# the p-values of the tests:
pShap

##          ext          ori
## 0.6650076 0.5384419
```

These results indicate that at the stage-level resolution and using the per capita rates (Foote 1999), both origination and extinction rate series have lognormal distributions, suggesting that mass extinction intervals are not qualitatively different from background turnover.

D. Which time slices are outliers from the distribution?

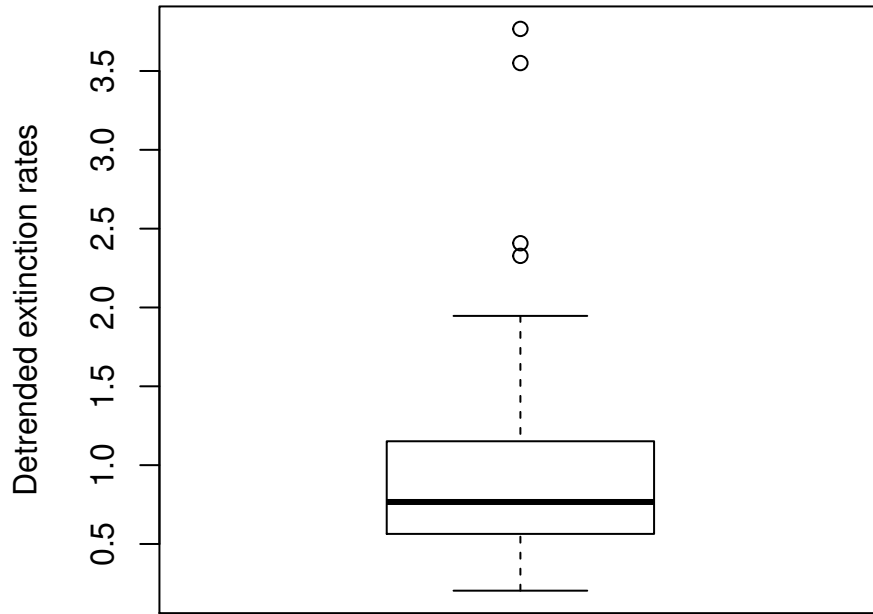
The exact definition of a ‘mass extinction interval’ has become somewhat blurry over the years. As rates fit a lognormal distribution quite well, it appears likely that background and mass extinction intervals are similar in qualitative terms, they are just at different positions along a spectrum. Even if this is the case, cataclysmic events (i.e. mass extinctions) will fall in the upper tail of the extinction rate distribution, which we can separate with non-parametric methods. The literature generally agrees that the largest ‘mass extinction’ was at the Permian/Triassic boundary, which we can reassess with a following line of code.

```
largest <- name[which(max(detExtShort, na.rm=T)==detExtShort)]
largest

## [1] "Changhsingian"
```

Which is the latest Permian interval, and is indistinguishable from the end-Permian mass extinction. We can also identify and contrast potential outliers with R’s `boxplot()` function.

```
boxp <- boxplot(detExtShort, ylab="Detrended extinction rates")
```

```
outliers <- name[as.numeric(names(boxp$out))]
outliers
```

```
## [1] "Katian"          "Changhsingian" "Rhaetian"       "Maastrichtian"
```

This suggests that without making any assumptions about the extinction rate distribution, four mass extinctions can be identified in the raw, stage-resolution data and the per capita rates. Among these, the end-Permian (Changhsingian), end-Triassic (Rhaetian) and end-Cretaceous (Maastrichtian) mass extinctions ('Big Three' in Alroy, 2008) are the most studied and usually represent the highest values in detrended extinction rate series. The presences of these three events are assessed with different rate calculation methods using the function presented in section 5.3. The Katian value is most likely associated with the end-Ordovician event.

E. Can we find traces of equilibrial dynamics in the series?

The quest to detect equilibrial patterns of diversity dynamics is rooted in island biogeography. Equilibrial dynamics means that diversity is bounded, but it does not necessarily mean that carrying capacity is stable through time. It rather points to the fact that there is an attractor in the richness dimension, pulling diversity up, when it is relatively low, and depressing it when it is relatively high (Alroy, 2010b). Alroy (2008) argued that equilibrial dynamics should manifest in three cross-correlations:

- Higher origination rates in time slice i should lead to higher diversity in time slice $i+1$.
- Higher diversities in time slice i in general should lead to higher extinction rates in time slice $i+1$.

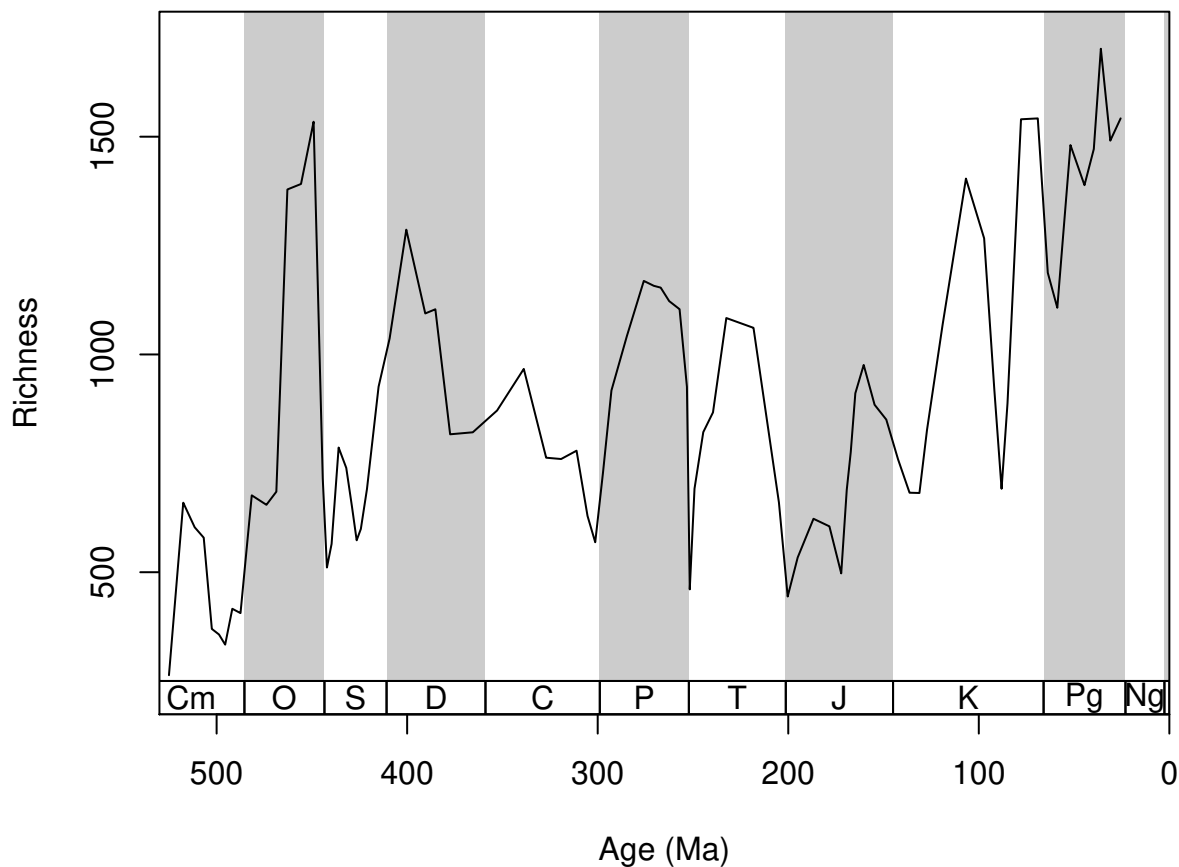
- Higher extinction rates in time slice i should lead to higher origination rate in slice $i+1$.

We have already detrended the turnover rates, but we still need to detrend the richness values. We follow the same basic procedure.

```
# select the same interval as for turnover rates
divVarLong <- div
divVarLong[indNAlong] <- NA
```

This can be plotted with:

```
# plotting the logged series
nMaxDiv <- max(divVarLong, na.rm=T)
nMinDiv <- min(divVarLong, na.rm=T)
tsplot(stages, boxes="per", shading="per", ylim=c(nMinDiv*0.95, nMaxDiv*1.05),
       ylab="Richness", xlim=c(530,0), xlab="Age (Ma)")
lines(age, divVarLong, col="black")
```



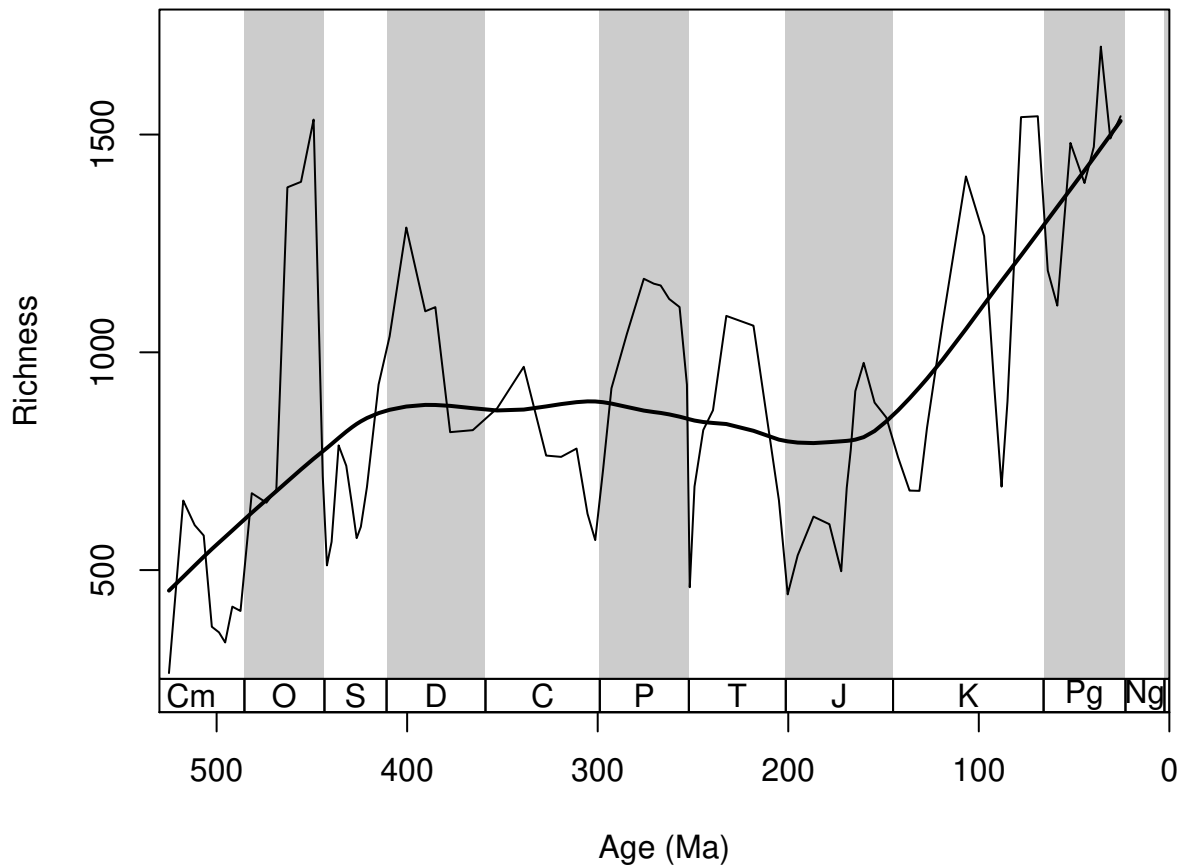
After this step, the LOESS regression is applied to the data.

```
divMiss <- !is.na(divVarLong)
transDivNoNA <- divVarLong[divMiss]
ageDiv <- age[divMiss]
```

```
divModel <- fANCOVA::loess.as(ageDiv, transDivNoNA, degree = 1,
  criterion = "aicc", user.span = NULL, plot = FALSE)
```

The predictions of this model are then calculated.

```
# predicted
transPredict <- predict(divModel, newdata=data.frame(x=age))
lines(age, transPredict, lwd=2)
```



The residuals are taken and are rescaled to the original magnitude, using the mean of the original series.

```
transResid <- divVarLong/transPredict # multiplicative decomposition
detDiv <- mean(divVarLong, na.rm=T)*transResid # rescaling
```

The series was then subsetting to include only those values that are in the detrended turnover rate series.

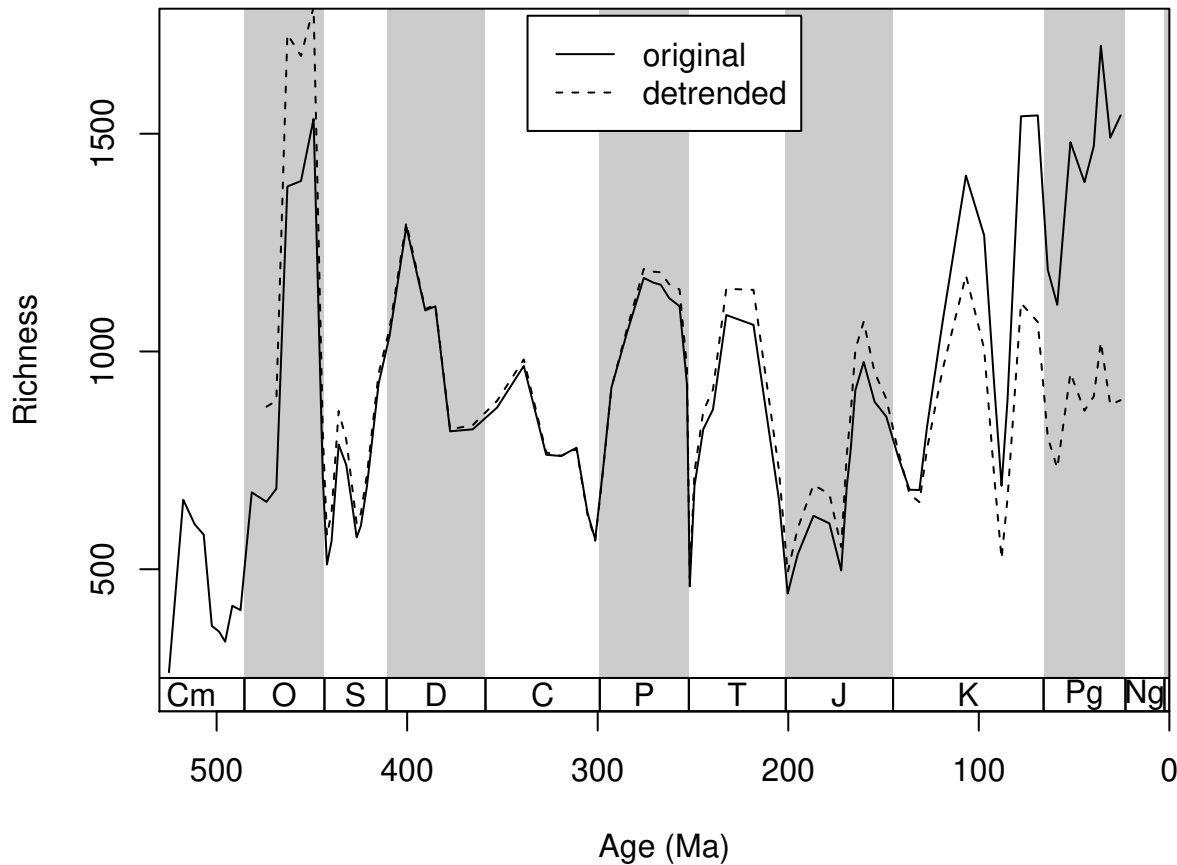
```
# limit to the same part as the turnover
detDivShort <- detDiv
detDivShort[indNashort] <- NA
```

The resulting series is comparable to the original series, but without the increasing trajectory.

```

# y ranges
tsplot(stages, boxes="per", shading="per", ylim=c(nMinDiv*0.95, nMaxDiv*1.05),
       ylab="Richness", xlim=c(530,0), xlab="Age (Ma)")
lines(age, divVarLong, lty=1)
lines(age, detDivShort, col="black", lty=2)
legend("top", inset=c(0.01, 0.01), legend=c("original", "detrended"), lty=c(1,2),
      bg="white")

```



After the diversity series is detrended, another small function `dynamics()` is applied to calculate cross correlation patterns at the chosen lag between the detrended extinction rate, origination rate and diversity variables. You can find this in the `phanDyn.R` file. `ext`, `ori` and `div` are the original detrended series (extinction, origination and diversity, respectively), and `extPlus`, `oriPlus` and `divPlus` denote the shifted series. For instance, a positive correlation between `ext` and `oriPlus` indicates that high extinction rate values are usually followed by high origination rate values in the next bin. Correlations between `ext` and `extPlus` indicate autocorrelations. A single run of this function produces the following results, with the alpha level of 0.01.

```

dyn <- dynamics(ori=detOriShort, ext=detExtShort, div=detDivShort, method="spearman",
               l=1, alpha=0.01)
# the $sig element returns significant components

```

```
dyn$sig
```

```
##      Var1    Var2 estimate      p
## 2      div divPlus 0.6488412 0.000000e+00
## 3      div     ori 0.4171550 2.236472e-04
## 5      div     ext 0.6255192 2.845641e-09
## 6      div extPlus 0.6048871 2.312953e-08
## 9  divPlus     ori 0.3872788 7.131964e-04
## 10 divPlus oriPlus 0.4171550 2.236472e-04
## 12 divPlus extPlus 0.6255192 2.845641e-09
## 16     ori oriPlus 0.3580155 1.846502e-03
## 30     ext extPlus 0.4078934 3.466648e-04
```

5.3. Single analysis function

After the basic analytical script was completed to answer the questions in sections 5.1 and 5.2 (above) a single function was prepared that takes any origination rate, extinction rate and richness time series and performs the analyses that we presented above. This function `analyzeMetrics()` can be found in the `phanDyn.R` file. The function requires a single `data.frame` object that has to include variables with the three time series, the bin durations, bin midpoint ages and names. All this information can be compiled by concatenating the time scale object and the output `data.frame` of the `divDyn()` function.

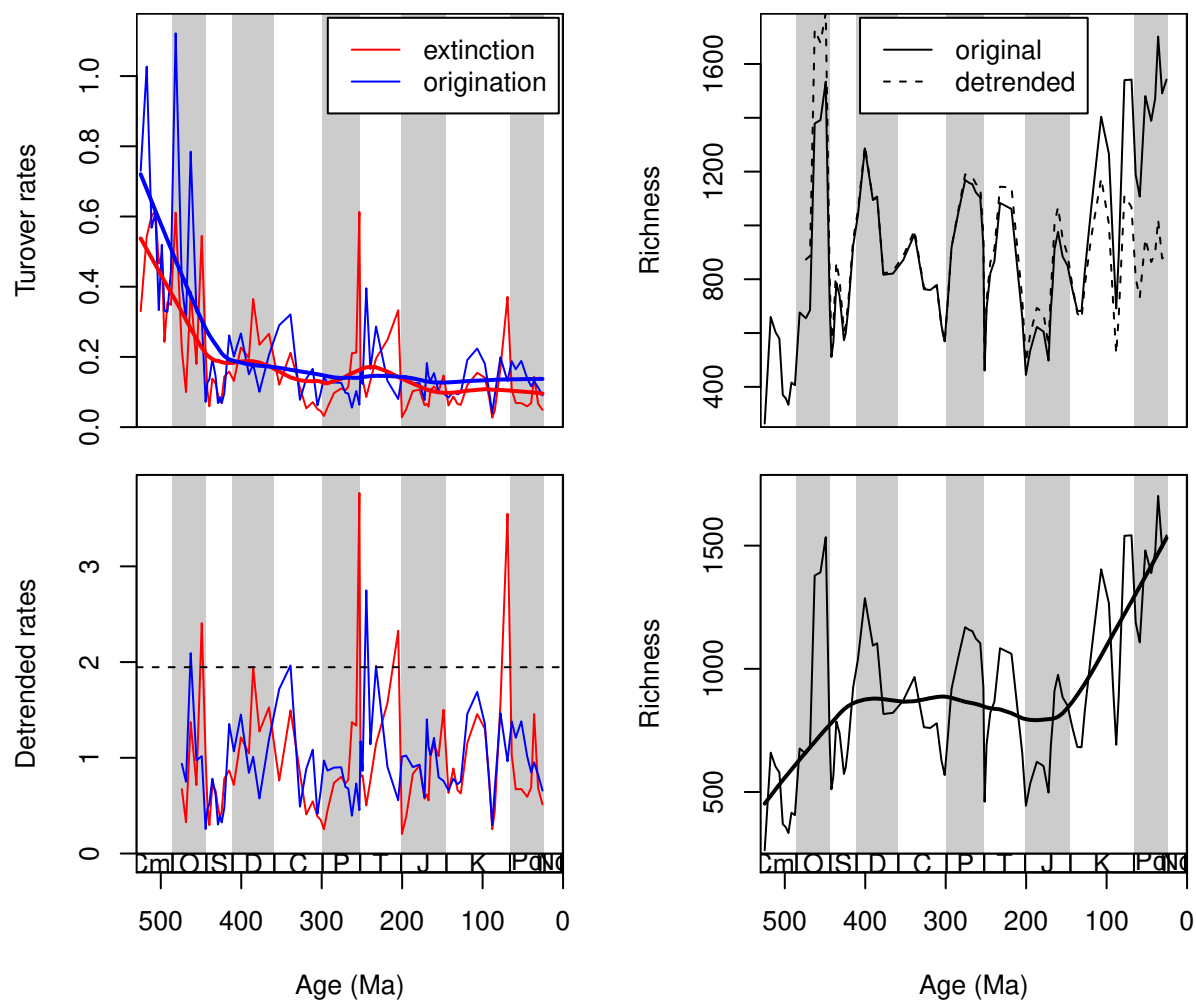
```
metrics <- cbind(ddStages, stages)
```

The arguments `ext` (extinction), `ori` (origination), `div` (richness/diversity), `age`, `dur` (duration) and `name` refer to the column names of the required variables. The function was implemented to return whether the ‘pulsed’ or ‘continuous’ model is supported, and to apply normalization with bin durations, if it is chosen so by the user. Normalized rates are only used in further analyses, when a correlation between rates and bin durations was evident in the non-normalized data and when the argument `normalize` is set to `TRUE`. In case both correlation tests are significant, the pulsed model is used as the default. Setting this argument to `FALSE` will coerce the calculations to use the input time series, using the ‘pulsed’ model throughout.

As mentioned in section 5.2, there are multiple ways to detrend the series. We implemented some of these methods, the procedure can be configured by adjusting the `detrend`, `transform` and `additive` arguments of the `analyzeMetrics()` function. The `detrend` argument specifies a switch between the different detrending options. Setting this to `"loess"` will run the detrending process shown in section 5.2 based on LOESS regression with the `fANCOVA` package (Wang, 2010). The option `"linear"` will fit a linear model to all three series. Setting the argument to `"arima"` will install and use the `forecast` package (Hyndman and Khandakar, 2008) to fit ARIMA models with the `auto.arima()` function. The argument `transform` specifies which transformation should be applied to the time series before the decomposition takes place. Setting this argument to `FALSE` will use the original series, `"log"` applies logarithm transformation and then exponentiation, `"sqrt"` applies square root transformation and then the squaring of the predicted and detrended values. The logical argument `additive` indicates whether additive or multiplicative decomposition is used to remove the trend.

The function also has additional arguments that enable the plotting of time series (`plot=TRUE`) and the output of messages and warnings (`feedback=TRUE`). If `plot` is set to `TRUE` a single four-panel plot is produced, where the panels depict: 1. rate values, 2. richness values, 3. detrended rates, 4. transformed richness values (identical to 3., if no transformation happens).

```
res <- analyzeMetrics(metrics, ext="extPC", ori="oriPC", div="divCSIB",
  age="mid", dur="dur", name="name", normalize=FALSE,
  plot=TRUE, feedback=FALSE)
```



The list class output of the function includes the objects that were introduced earlier in sections 5.1 and 5.2.

```
res
```

```
## $`whichModel` is indicated, pulsed or continuous`
## [1] TRUE
##
## $`pulsed-continuous significance`
##
##               est                p
## extinction (not-normalized) 0.2311910 2.835024e-02
## extinction (normalized)    -0.4545859 6.767625e-06
## origination (not-normalized) 0.3462226 8.298011e-04
## origination (normalized)    -0.5206346 1.439519e-07
##
## $declines
##
##               correlation          p-value
## extinctions          0.52126246 4.794109e-07
## post-Ordovician extinctions 0.23094629 5.642085e-02
```

```

## originations                0.44487004 2.435087e-05
## post-Ordovician originations 0.06024845 6.221282e-01
##
## `$largest extinction`
## [1] "end-Permian"
##
## `$extinction outliers (boxplot)`
## [1] "other"          "end-Permian"      "end-Triassic"     "end-Cretaceous"
##
## `$log-normality (Shapiro-Wilk p-value)`
##      ext      ori
## 0.6650076 0.5384419
##
## $dyn
## $dyn$est
##      div      divPlus      ori      oriPlus      ext
## div      1.000000000 0.6488412 0.4171550 -0.001377268 0.625519203
## divPlus  0.648841170 1.0000000 0.3872788  0.417155050 0.108419104
## ori      0.417155050 0.3872788 1.0000000  0.358015550 0.295078236
## oriPlus -0.001377268 0.4171550 0.3580155  1.000000000 0.005997779
## ext      0.625519203 0.1084191 0.2950782  0.005997779 1.000000000
## extPlus  0.604887079 0.6255192 0.2965864  0.295078236 0.407893373
##      extPlus
## div      0.6048871
## divPlus  0.6255192
## ori      0.2965864
## oriPlus  0.2950782
## ext      0.4078934
## extPlus  1.0000000
##
## $dyn$p
##      div      divPlus      ori      oriPlus      ext
## div      0.000000e+00 0.000000e+00 0.0002236472 0.9907847779 2.845641e-09
## divPlus  0.000000e+00 0.000000e+00 0.0007131964 0.0002236472 3.570489e-01
## ori      2.236472e-04 7.131964e-04 0.0000000000 0.0018465021 1.040622e-02
## oriPlus  9.907848e-01 2.236472e-04 0.0018465021 0.0000000000 9.595493e-01
## ext      2.845641e-09 3.570489e-01 0.0104062175 0.9595493415 0.000000e+00
## extPlus  2.312953e-08 2.845641e-09 0.0105321510 0.0104062175 3.466648e-04
##      extPlus
## div      2.312953e-08
## divPlus  2.845641e-09
## ori      1.053215e-02
## oriPlus  1.040622e-02
## ext      3.466648e-04
## extPlus  0.000000e+00
##
## $dyn$sig
##      Var1      Var2 estimate      p
## 2      div divPlus 0.6488412 0.000000e+00
## 3      div      ori 0.4171550 2.236472e-04
## 5      div      ext 0.6255192 2.845641e-09
## 6      div extPlus 0.6048871 2.312953e-08
## 9  divPlus      ori 0.3872788 7.131964e-04
## 10 divPlus oriPlus 0.4171550 2.236472e-04

```

```
## 12 divPlus extPlus 0.6255192 2.845641e-09
## 16      ori oriPlus 0.3580155 1.846502e-03
## 17      ext      ori 0.2950782 1.040622e-02
## 18 extPlus      ori 0.2965864 1.053215e-02
## 24 extPlus oriPlus 0.2950782 1.040622e-02
## 30      ext extPlus 0.4078934 3.466648e-04
```

6. Applying the analyses to multiple series

The best thing about the rate and diversity estimators is also the worst thing: they are in continuous development. Although there is evidence to support the better applicability of some methods (Alroy, 2014; 2015), different researchers may favor different solutions to answer scientific hypotheses. Therefore, we are not trying to impose the use of any methods, but rather intend to make these calculations available to everybody and demonstrate the effects of methodological choices on the results.

For instance, results depend on the chosen equations to calculate the taxonomic turnover series. Although the per capita (Foote, 1999) rates were the most frequently used method in the past two decades, using these equations for the estimation of the rates can lead to systematic errors close to the edges of the time series (i.e. edge effects; Foote, 2000). Other methods (Alroy, 2008; 2014; 2015) were developed to be unaffected by this phenomenon. Nevertheless, even within the series, the different methods will lead to considerably varying values that can have an effect on the results used to test overarching hypotheses.

Therefore, we decided to redo the analysis presented in section 5 using different metrics of turnover and methods of sampling standardization, to assess how robust the general findings are in the face of different methodologies. Sets of time series were drafted based on two different resolutions (10my bin and stages), three different data treatments (using raw data, Classical Rarefaction and Shareholder Quorum Subsampling) and four different rate calculation methods (per capita rates, corrected three-timer rates, gap-filler equations and second-for-third substitution rates). Metrics published earlier were not considered in the analyses as they have known problems with accuracy (Foote, 1994). Equilibrial dynamics were tested using the corrected SIB diversity values.

6.1. Calculating time series with multiple methods

The candidate metrics are calculated with the different resolutions and data treatment options. First, the calculations with the 10 my time scale are implemented.

Before the calculation of the actual values can take place, the desired level of sampling intensities has to be determined for the sampling standardization processes. Classical Rarefaction (CR) is the oldest of the subsampling methods (Raup, 1975). Diversity curves drafted with this method are strongly dependent on sampling intensity: Alroy (2010b) has shown that the method leads to richness curves that flatten as the quota (desired sampling intensity in occurrences) decreases. Despite the fact that CR has its own limitations for richness estimation, it is still one of the most widely used sampling standardization protocol in the literature, and it produces comparable results to SQS with our data.

The quota for CR is set so that a complete time series can be produced without creating time slices with ‘failed’ subsampling, where there are not enough sampled occurrences. These were already calculated and are present in the `sampBins` and `sampStages` objects.

Let’s first look at the number of occurrences in the 10 myr stage-results. We need to find out which bins are the least sampled to figure out a good subsampling quota. To do this, it is very useful to assign the bin identifiers to the number of occurrences first (names attribute), and then we can put the occurrences in ascending order.


```
binOccs <- sampBin$occs
names(binOccs) <- paste("bin", rownames(sampBin), sep="")
# in ascending order
sort(binOccs)
```

```
## bin1 bin16 bin35 bin17 bin5 bin38 bin18 bin13 bin37 bin36 bin3 bin30
## 568 4831 4852 4859 5247 5317 5820 6411 6439 6601 6626 6698
## bin28 bin21 bin15 bin19 bin39 bin20 bin4 bin45 bin40 bin6 bin44 bin7
## 6793 6808 7806 7882 8449 9114 9158 9187 9895 10234 10467 10496
## bin33 bin12 bin42 bin2 bin41 bin29 bin31 bin22 bin47 bin26 bin27 bin32
## 10637 10740 10766 11172 11484 11684 12180 12185 12364 12601 12693 12699
## bin46 bin10 bin9 bin23 bin34 bin25 bin48 bin11 bin14 bin43 bin24 bin8
## 14677 14768 16013 20049 23197 24504 24533 24665 25348 26906 29042 32830
## bin49
## 49772
```

This indicates that the first bin is disproportionately poorly sampled. This is unlikely to convey important information for most of the time series. At the same time, subsampling the whole series to this low sampling level would destroy a huge amount of information. However, the second lowest level of about 4,800 occurrences (bin 16, Devonian 5) looks like a decent level of sampling and can be applied to almost the whole time series.

The same reasoning can be applied to the stage-level data.

```
stgOccs <- sampStg$occs
names(stgOccs) <- paste("stg", rownames(sampStg), sep="")
sort(stgOccs)
```

```
## stg4 stg11 stg9 stg78 stg5 stg13 stg26 stg59 stg41 stg95 stg70 stg7
## 1154 1444 1474 1595 1636 1767 1962 2221 2314 2505 2546 2555
## stg22 stg63 stg73 stg39 stg10 stg20 stg28 stg25 stg12 stg27 stg72 stg38
## 2654 2911 2991 3024 3055 3096 3114 3167 3303 3311 3319 3499
## stg16 stg83 stg79 stg55 stg85 stg60 stg71 stg21 stg30 stg8 stg42 stg15
## 3514 3566 3590 3676 3712 3830 3857 3874 3993 4058 4205 4244
## stg64 stg58 stg89 stg23 stg43 stg35 stg69 stg36 stg82 stg14 stg52 stg68
## 4279 4523 4639 4715 4759 4831 4852 4859 4929 4993 5126 5142
## stg6 stg84 stg77 stg40 stg74 stg32 stg65 stg24 stg57 stg86 stg29 stg88
## 5220 5260 5297 5314 5317 5370 5475 5985 6053 6211 6227 6396
## stg31 stg53 stg56 stg44 stg45 stg66 stg34 stg67 stg87 stg75 stg54 stg48
## 6411 6785 6793 6821 7359 7586 7806 7998 8055 8449 8619 8971
## stg49 stg47 stg62 stg76 stg80 stg51 stg91 stg90 stg37 stg61 stg46 stg50
## 9074 9149 9583 9895 10766 10932 11083 11185 11243 12180 12261 12306
## stg18 stg17 stg92 stg93 stg94 stg33 stg81 stg19
## 12348 12552 15463 15893 16143 17847 26906 27644
```

In this case, the number of occurrences in a bin increases gradually, and multiple poorly sampled bins are in the middle of the series. Using different quotas would create different results, for the sake of simplicity, it is probably best to restrain ourselves to a lower quota of 1,100 occurrences for the stage-level analyses.

For SQS, the subsampling configuration discussed in Section 4 stands here, too. A major advantage of SQS to CR is that it gives you more or less the same curve regardless of the quorum, as long as the quorum isn't extremely low (< 0.4 or so, Alroy 2010b).

The large quantity of data and the relatively small subsample sizes indicate that one must carry out a high number of iterations for the estimates to stabilize. The results presented in this section are based on 300 iterations, which will take a considerable amount of time to run, but you can decrease this parameter at will by setting the `iter` argument of the `subsample()` function.

```

# 1. raw patterns
ddBins <- divDyn(dat, bin="bin", tax="clgen")
# 2. CR
crBins <- subsample(dat, bin="bin", tax="clgen", coll="collection_no", q=4800, iter=300,
  duplicates=FALSE)
# 3. SQS
sqsBins <- subsample(dat, bin="bin", tax="clgen", coll="collection_no", q=0.7, iter=300,
  ref="reference_no", singleton="ref", type="sqs", duplicates=FALSE, excludeDominant=TRUE,
  largestColl =TRUE)

```

Then the calculations were repeated for the stage-level resolution.

```

# 1. raw patterns
ddStg <- divDyn(dat, bin="stg", tax="clgen")
# 2. CR
crStg <- subsample(dat, bin="stg", tax="clgen", coll="collection_no", q=1100,
  iter=300, duplicates=FALSE)
# 3. SQS
sqsStg <- subsample(dat, bin="stg", tax="clgen", coll="collection_no", q=0.7,
  iter=300, ref="reference_no", singleton="ref", type="sqs", duplicates=FALSE,
  excludeDominant=TRUE, largestColl =TRUE)

```

In total, we have 6 different result objects.

6.2. Running the analysis script on the different sets of time series

The next thing is to organize the results (i.e., create identifiers of the different time series) so the analyses described in section 5 can be iterated. Each run of the analytical script needs three time series (origination, extinction, richness), which can be identified with the appropriate result table names (stratigraphic resolution, data treatment) and the variable names (metric type). The following object `comb` is created to organize these combinations.

```

# combination table
# types of rates
comb <- matrix(
  c(
    "extPC", "oriPC", "divCSIB", # per capita rates
    "extC3t", "oriC3t", "divCSIB", # corrected 3t rates
    "extGF", "oriGF", "divCSIB", # gap-filler rates
    "ext2f3", "ori2f3", "divCSIB" # second-for-third-substitution rates
  ),
  ncol=3, nrow=4, byrow=T)
rownames(comb) <- c("PC", "C3t", "GF", "2f3")
comb <- comb[rep(1:4, 6), ]

# the result matrices
sourceVar<-rep(c("ddBins", "crBins", "sqsBins", "ddStg", "crStg", "sqsStg"),each=4)

# timescale objects
scale <- rep(c("bins", "stages"), each=12)
# combine everything together
comb <- cbind(sourceVar, scale, comb)
colnames(comb) <- c("source", "timescale", "ext", "ori", "div")

```

```
# the names of the results (rownames)
subtype <- rep(c("raw", "cr", "sqs"), each=4),2)
rownames(comb) <- paste(subtype, rownames(comb), sep="")
rownames(comb) <- paste(rownames(comb), rep(c("10my", "stages"), each=12), sep="_")
comb
```

```
##      source    timescale ext      ori      div
## rawPC_10my    "ddBins"   "bins"   "extPC" "oriPC" "divCSIB"
## rawC3t_10my   "ddBins"   "bins"   "extC3t" "oriC3t" "divCSIB"
## rawGF_10my    "ddBins"   "bins"   "extGF" "oriGF" "divCSIB"
## raw2f3_10my   "ddBins"   "bins"   "ext2f3" "ori2f3" "divCSIB"
## crPC_10my     "crBins"   "bins"   "extPC" "oriPC" "divCSIB"
## crC3t_10my    "crBins"   "bins"   "extC3t" "oriC3t" "divCSIB"
## crGF_10my     "crBins"   "bins"   "extGF" "oriGF" "divCSIB"
## cr2f3_10my    "crBins"   "bins"   "ext2f3" "ori2f3" "divCSIB"
## sqsPC_10my    "sqsBins"  "bins"   "extPC" "oriPC" "divCSIB"
## sqsC3t_10my   "sqsBins"  "bins"   "extC3t" "oriC3t" "divCSIB"
## sqsGF_10my    "sqsBins"  "bins"   "extGF" "oriGF" "divCSIB"
## sqs2f3_10my   "sqsBins"  "bins"   "ext2f3" "ori2f3" "divCSIB"
## rawPC_stages  "ddStg"    "stages" "extPC" "oriPC" "divCSIB"
## rawC3t_stages "ddStg"    "stages" "extC3t" "oriC3t" "divCSIB"
## rawGF_stages  "ddStg"    "stages" "extGF" "oriGF" "divCSIB"
## raw2f3_stages "ddStg"    "stages" "ext2f3" "ori2f3" "divCSIB"
## crPC_stages   "crStg"    "stages" "extPC" "oriPC" "divCSIB"
## crC3t_stages  "crStg"    "stages" "extC3t" "oriC3t" "divCSIB"
## crGF_stages   "crStg"    "stages" "extGF" "oriGF" "divCSIB"
## cr2f3_stages  "crStg"    "stages" "ext2f3" "ori2f3" "divCSIB"
## sqsPC_stages  "sqsStg"   "stages" "extPC" "oriPC" "divCSIB"
## sqsC3t_stages "sqsStg"   "stages" "extC3t" "oriC3t" "divCSIB"
## sqsGF_stages  "sqsStg"   "stages" "extGF" "oriGF" "divCSIB"
## sqs2f3_stages "sqsStg"   "stages" "ext2f3" "ori2f3" "divCSIB"
```

In this table, every row corresponds to the arguments of a certain set of time series. For instance, the first row points to the raw per capita rates with the 10 myr stratigraphic resolution. Now, all that remain is to initialize and iterate the analyses using the different sets of time series. In order to display the rates together, they have to be saved in separate containers.

```
# matrices to hold rates for later plotting
extDatBin <- data.frame(bins=1:49) # extinctions, 10my
oriDatBin <- data.frame(bins=1:49) # originations, 10my
extDatStg <- data.frame(stg=1:95) # extinctions, stages
oriDatStg <- data.frame(stg=1:95) # originations, stages
```

The first two `data.frame` objects correspond to the 10my bin resolution and the second one to the stage-level resolution. The following code applies the function to the different time series in a `for` loop (for transparency).

```
#! # initialize pdf, if you want to
#! pdf("2018-08-21_marineAnimals2_redo.PDF", 17,14)
# iterate through all cases
for(i in 1:nrow(comb)){
  # name of the set
  case <- rownames(comb)[i]
  # the results matrix and the timescale object combined
  metrics <- cbind(get(comb[i, "timescale"]), get(comb[i, "source"]))

  # save rates for later, depending on the resolution
```

```

# 10 my
if(comb[i, "timescale"] == "bins"){
  oriDatBin[[case]] <- metrics[, comb[i, "ori"]]
  extDatBin[[case]] <- metrics[, comb[i, "ext"]]
}
# stages
if(comb[i, "timescale"] == "stages"){
  oriDatStg[[case]] <- metrics[, comb[i, "ori"]]
  extDatStg[[case]] <- metrics[, comb[i, "ext"]]
}
# the analytical function
res <- analyzeMetrics(metrics, ext=comb[i, "ext"], ori=comb[i, "ori"],
  div=comb[i, "div"], age="mid", dur="dur", name="name", normalize=TRUE,
  plot=FALSE, feedback=FALSE, detrend="loess", transform=FALSE, additive=FALSE)

# save in global namespace with unique name
assign(case, res)

# add the name to the plot
#! par(mfrow=c(1,1))
#! mtext(side=3, text=case, line=-2, cex=3)
}

#! dev.off()

```

With this snippet, each result will be stored as a list class object in the global namespace, with the name of the row in `comb` that contains the arguments.

All lines commented with `#!` are part of the embedded plotting functionality. If you uncomment these lines (delete `!` too!) and set `plot=TRUE` in the `analyzeMetrics()` function call, a single .pdf file will be produced that will show the four-panel figure presented in section 5.3 for every set of time series (row in `comb`) on a separate page. You can take look at this .pdf file by following this link:

<https://github.com/adamkocsis/ddPhanero/raw/master/export/0.3/detrending.pdf>

6.3. Summarizing the results (Table 2)

To get an idea of the generality of the results, the support for the individual hypotheses must be visualized in a single display item. We prepared a function to extract correlation coefficient estimates and p -values from the individual objects. Two temporary tables were prepared – one for the estimates and another one for the significance values. The idea behind this approach is to use the first table to present the actual values, and use the table of p -values to format the first table (using conditional formatting in MS Excel). The function to gather the relevant information from the lists can be inspected in the `phanDyn.R` file (`extractVals()` function).

```

# the shown values
values <- extractVals(rownames(comb), pvals=FALSE)

```

This is a fairly large table that includes every case-related results in a column. For correlation tests, the table denotes the coefficient estimates. For the presence of mass-extinction intervals, the table shows binary values. As the iteration through the output object is the same for the p -values and the estimates, it was a straightforward choice to use the same function to extract the p -values. Setting the `pvals` argument of this function to `TRUE` will extract the p -values, where another they are also present.

```
# the p values
pvals <- extractVals(rownames(comb), pvals=TRUE)
```

Note that the entries for the presences of mass extinctions in this table there are no significance values. The table was compiled to be used for conditional formatting and 0.002 values indicate the presence of a mass extinction, 0.02 indicates its absence.

The support or rejection of hypotheses depend on the p -values. The results based on the 10 myr-scale results have the following values (rounded and transposed for better readability):

```
round(t(pvals[1:12,]),3)
```

##	rawPC_10my	rawC3t_10my	rawGF_10my
## ext. rates with durations	0.671	0.805	0.670
## norm. ext. rates with durations	0.000	0.002	0.018
## orig. rates with durations	0.099	0.165	0.026
## norm. orig. rates with durations	0.011	0.004	0.182
## extinctions	0.000	0.000	0.000
## post-Ordovician extinctions	0.063	0.019	0.009
## originations	0.001	0.002	0.002
## post-Ordovician originations	0.153	0.146	0.764
## end-Permian ME	0.002	0.002	0.002
## end-Triassic ME	0.002	0.002	0.002
## end-Cretaceous ME	0.002	0.002	0.002
## end-Permian is highest	0.002	0.020	0.002
## number of mass extinctions	0.002	0.002	0.002
## extinctions log-normal (p-values)	0.139	0.384	0.053
## originations log-normal (p-values)	0.234	0.551	0.883
## origination and lagged diversity	0.067	0.388	0.285
## diversity and lagged extinction	0.003	0.065	0.117
## diversity and lagged origination	0.811	0.000	0.084
##	raw2f3_10my	crPC_10my	crC3t_10my
## ext. rates with durations	0.999	0.694	0.302
## norm. ext. rates with durations	0.041	0.000	0.000
## orig. rates with durations	0.110	0.129	0.454
## norm. orig. rates with durations	0.207	0.015	0.006
## extinctions	0.000	0.000	0.000
## post-Ordovician extinctions	0.139	0.078	0.004
## originations	0.063	0.000	0.004
## post-Ordovician originations	0.169	0.066	0.108
## end-Permian ME	0.002	0.002	0.002
## end-Triassic ME	0.002	0.002	0.002
## end-Cretaceous ME	0.002	0.020	0.002
## end-Permian is highest	0.002	0.002	0.002
## number of mass extinctions	0.002	0.002	0.002
## extinctions log-normal (p-values)	0.684	0.080	0.045
## originations log-normal (p-values)	0.001	0.830	0.923
## origination and lagged diversity	0.070	0.028	0.180
## diversity and lagged extinction	0.111	0.026	0.298
## diversity and lagged origination	0.533	0.488	0.000
##	crGF_10my	cr2f3_10my	sqsPC_10my
## ext. rates with durations	0.635	0.394	0.977
## norm. ext. rates with durations	0.000	0.000	0.001
## orig. rates with durations	0.214	0.505	0.075
## norm. orig. rates with durations	0.101	0.064	0.065

## extinctions	0.000	0.000	0.000
## post-Ordovician extinctions	0.008	0.027	0.082
## originations	0.000	0.005	0.000
## post-Ordovician originations	0.067	0.064	0.128
## end-Permian ME	0.002	0.002	0.002
## end-Triassic ME	0.002	0.002	0.002
## end-Cretaceous ME	0.002	0.002	0.020
## end-Permian is highest	0.002	0.002	0.002
## number of mass extinctions	0.002	0.002	0.002
## extinctions log-normal (p-values)	0.078	0.081	0.474
## originations log-normal (p-values)	0.994	0.350	0.791
## origination and lagged diversity	0.083	0.043	0.003
## diversity and lagged extinction	0.139	0.165	0.002
## diversity and lagged origination	0.003	0.006	0.348
##	sqsC3t_10my	sqsGF_10my	sqs2f3_10my
## ext. rates with durations	0.500	0.518	0.989
## norm. ext. rates with durations	0.001	0.023	0.018
## orig. rates with durations	0.281	0.084	0.288
## norm. orig. rates with durations	0.014	0.383	0.315
## extinctions	0.000	0.000	0.000
## post-Ordovician extinctions	0.015	0.005	0.063
## originations	0.003	0.000	0.057
## post-Ordovician originations	0.211	0.140	0.142
## end-Permian ME	0.002	0.002	0.002
## end-Triassic ME	0.002	0.002	0.002
## end-Cretaceous ME	0.002	0.020	0.020
## end-Permian is highest	0.002	0.002	0.002
## number of mass extinctions	0.002	0.002	0.002
## extinctions log-normal (p-values)	0.298	0.378	0.613
## originations log-normal (p-values)	0.902	0.023	0.012
## origination and lagged diversity	0.054	0.051	0.021
## diversity and lagged extinction	0.187	0.099	0.090
## diversity and lagged origination	0.000	0.205	0.892

The stage-level *p*-values can also be extracted in a similar way:

```
round(t(pvals[13:24,]),3)
```

##	rawPC_stages	rawC3t_stages	rawGF_stages
## ext. rates with durations	0.028	0.175	0.036
## norm. ext. rates with durations	0.000	0.001	0.001
## orig. rates with durations	0.001	0.070	0.002
## norm. orig. rates with durations	0.000	0.000	0.000
## extinctions	0.000	0.000	0.000
## post-Ordovician extinctions	0.056	0.259	0.243
## originations	0.000	0.001	0.000
## post-Ordovician originations	0.622	0.681	0.663
## end-Permian ME	0.002	0.002	0.002
## end-Triassic ME	0.002	0.002	0.002
## end-Cretaceous ME	0.002	0.002	0.002
## end-Permian is highest	0.002	0.002	0.002
## number of mass extinctions	0.002	0.002	0.002
## extinctions log-normal (p-values)	0.665	0.171	0.012
## originations log-normal (p-values)	0.538	0.047	0.000
## origination and lagged diversity	0.001	0.011	0.073

## diversity and lagged extinction	0.000	0.000	0.001
## diversity and lagged origination	0.960	0.248	0.137
##	raw2f3_stages	crPC_stages	crC3t_stages
## ext. rates with durations	0.007	0.437	0.166
## norm. ext. rates with durations	0.106	0.000	0.000
## orig. rates with durations	0.001	0.098	0.137
## norm. orig. rates with durations	0.003	0.000	0.000
## extinctions	0.000	0.000	0.000
## post-Ordovician extinctions	0.037	0.026	0.297
## originations	0.001	0.000	0.000
## post-Ordovician originations	0.826	0.041	0.537
## end-Permian ME	0.002	0.002	0.002
## end-Triassic ME	0.002	0.002	0.002
## end-Cretaceous ME	0.002	0.002	0.002
## end-Permian is highest	0.002	0.002	0.002
## number of mass extinctions	0.002	0.002	0.002
## extinctions log-normal (p-values)	0.245	0.098	0.082
## originations log-normal (p-values)	0.357	0.957	0.671
## origination and lagged diversity	0.095	0.003	0.142
## diversity and lagged extinction	0.099	0.000	0.308
## diversity and lagged origination	0.810	0.222	0.000
##	crGF_stages	cr2f3_stages	sqsPC_stages
## ext. rates with durations	0.143	0.038	0.220
## norm. ext. rates with durations	0.000	0.000	0.000
## orig. rates with durations	0.154	0.202	0.028
## norm. orig. rates with durations	0.000	0.000	0.000
## extinctions	0.000	0.000	0.000
## post-Ordovician extinctions	0.365	0.250	0.015
## originations	0.000	0.000	0.000
## post-Ordovician originations	0.499	0.402	0.220
## end-Permian ME	0.002	0.002	0.002
## end-Triassic ME	0.002	0.002	0.002
## end-Cretaceous ME	0.002	0.002	0.002
## end-Permian is highest	0.002	0.002	0.002
## number of mass extinctions	0.002	0.002	0.002
## extinctions log-normal (p-values)	0.000	0.068	0.824
## originations log-normal (p-values)	0.006	0.001	0.580
## origination and lagged diversity	0.264	0.689	0.000
## diversity and lagged extinction	0.096	0.253	0.000
## diversity and lagged origination	0.000	0.000	0.967
##	sqsC3t_stages	sqsGF_stages	
## ext. rates with durations	0.191	0.051	
## norm. ext. rates with durations	0.000	0.000	
## orig. rates with durations	0.103	0.023	
## norm. orig. rates with durations	0.000	0.000	
## extinctions	0.000	0.000	
## post-Ordovician extinctions	0.411	0.458	
## originations	0.001	0.000	
## post-Ordovician originations	0.864	0.807	
## end-Permian ME	0.002	0.002	
## end-Triassic ME	0.002	0.002	
## end-Cretaceous ME	0.002	0.002	
## end-Permian is highest	0.002	0.002	
## number of mass extinctions	0.002	0.002	

## extinctions log-normal (p-values)	0.078	0.000
## originations log-normal (p-values)	0.360	0.001
## origination and lagged diversity	0.003	0.004
## diversity and lagged extinction	0.041	0.001
## diversity and lagged origination	0.216	0.388
##	sqs2f3_stages	
## ext. rates with durations	0.010	
## norm. ext. rates with durations	0.007	
## orig. rates with durations	0.034	
## norm. orig. rates with durations	0.000	
## extinctions	0.001	
## post-Ordovician extinctions	0.388	
## originations	0.000	
## post-Ordovician originations	0.628	
## end-Permian ME	0.002	
## end-Triassic ME	0.002	
## end-Cretaceous ME	0.002	
## end-Permian is highest	0.002	
## number of mass extinctions	0.002	
## extinctions log-normal (p-values)	0.003	
## originations log-normal (p-values)	0.023	
## origination and lagged diversity	0.018	
## diversity and lagged extinction	0.009	
## diversity and lagged origination	0.407	

The tables of *p*-values and estimates can be combined using conditional formatting to render a comprehensible table. The following formatting steps were applied to these tables that were later fed to an MS Excel spreadsheet.

```
# round most values
valuesRounded <- round(values,2)
# except for p-values
valuesRounded[,
  c("extinctions log-normal (p-values)",
    "originations log-normal (p-values)")] <-
  round(values[,c(
    "extinctions log-normal (p-values)",
    "originations log-normal (p-values)"]),3)

#transform to character
#valuesRounded <- as.data.frame(valuesRounded)
for(i in 9:13){
  valuesRounded[,i] <- as.character(valuesRounded[,i])
  valuesRounded[valuesRounded[,i]=="1",i] <- "yes"
  valuesRounded[valuesRounded[,i]=="0",i] <- "no"
}

# use inequality signs where values are very low
valuesRounded[
  valuesRounded[,"extinctions log-normal (p-values)"]=="0",
  "extinctions log-normal (p-values)"] <- "<0.001"
valuesRounded[
  valuesRounded[,"originations log-normal (p-values)"]=="0",
  "originations log-normal (p-values)"] <- "<0.001"
```


In the actual paper, Table 2 was compiled after the transposition of the estimate and *p*-value tables, and the separation of both tables to 10 myr time scale and stage-level timescale subsets.

```
# transpose everything
tValues<- t(valuesRounded)
tP <- t(pvals)

# separate based on resolution
# 10 myr
valBin <- tValues[,1:12]
pBin <- tP[,1:12]
# stages
valStage <- tValues[,13:24]
pStage <- tP[,13:24]
```

These four tables were then fed to the Excel spreadsheet we prepared to do the conditional formatting of the valBin and valStage tables. This MS Excel file (conditionalTable.xlsx) can also be found at the GitHub repository of the example.

6.4. Figure of rates with multiple methods (Figure 2)

To show the total variation in the taxonomic rate series, the origination and extinction rates were rendered at the two different resolutions. Four panels were drawn, each panel having 12 time series. The quantiles of the value distributions in the independent time slices were connected with the `shades()` function to provide backgrounds.

Then the panels in sequence:

```
par(mfrow=c(2,2))
# extinctions
# 1st panel - 10myr bins
par(mar=c(2, 5.1, 4.1, 1)) # margin adjustment
# plot
tsplot(stages, boxes="per", shading="per", ylim=c(0, 3), ylab="Extinction rates",
       xlim=4:95, plot.args=list(axes=F, cex.lab=0.8), xlab="", labels.args=list(cex=0.8))
axis(2, cex.axis=0.8)
shades(bins$mid, as.matrix(extDatBin[, -1]), col="red") # background
for(i in 2:ncol(extDatBin))
  lines(bins$mid, extDatBin[, i], lwd=0.5) # extinctions - BIN
mtext(line=1, text="10 my bins", side=3, cex=1) # show resolution

# 2nd panel - stages
par(mar=c(2, 1, 4.1, 4.1)) # margin adjustment
# plot
tsplot(stages, boxes="per", shading="per", ylim=c(0, 2), ylab="", xlim=4:95,
       plot.args=list(axes=F, cex.lab=0.8), xlab="", labels.args=list(cex=0.8))
shades(stages$mid, as.matrix(extDatStg[, -1]), col="red") # background
for(i in 2:ncol(extDatStg))
  lines(stages$mid, extDatStg[, i], lwd=0.5) # extinction - stages
mtext(line=1, text="stages", side=3, cex=1)

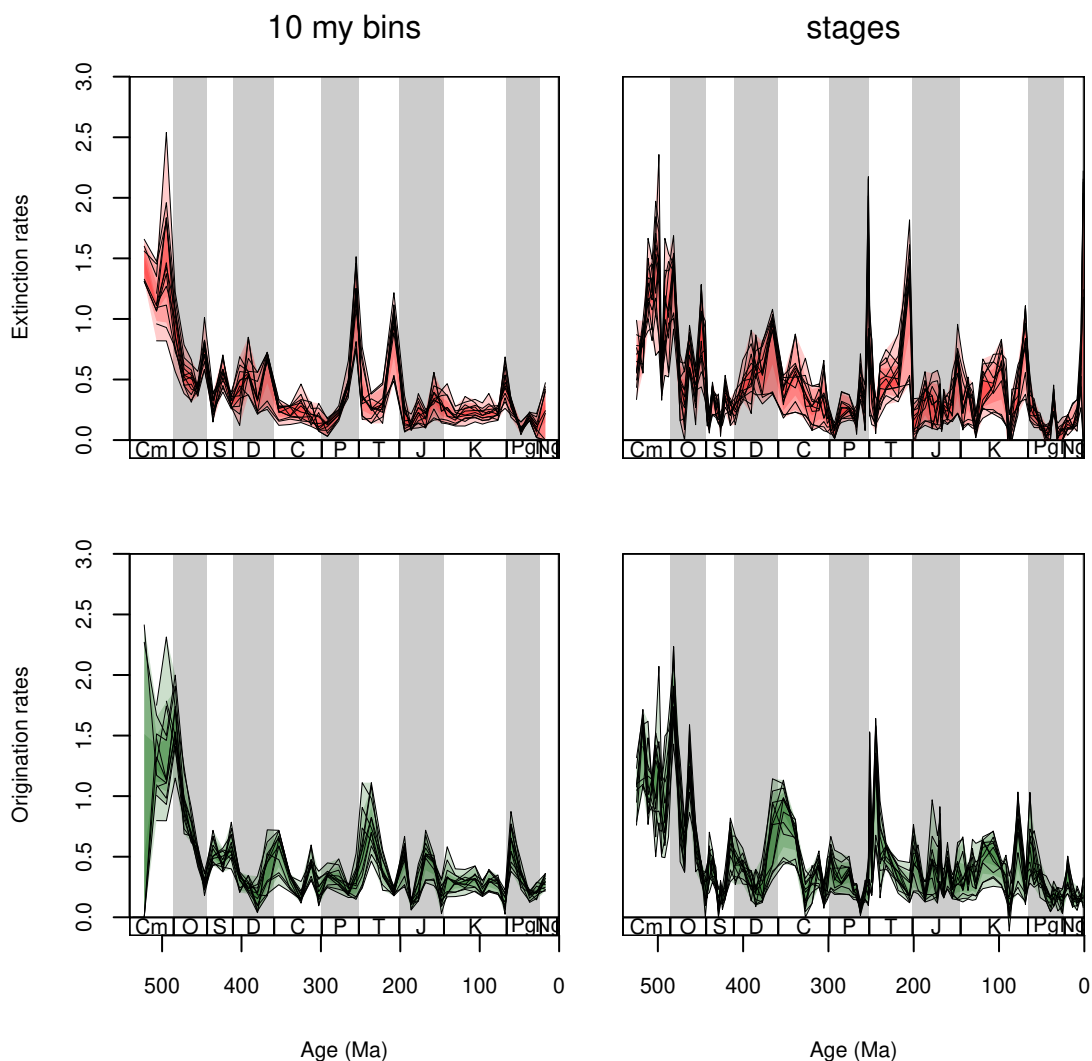
# originations
# 3rd panel - 10 myr bins
par(mar=c(5.1, 5.1, 1, 1)) # margin adjustment
# plot
```

```

tsplot(stages, boxes="per", shading="per", ylim=c(0, 3), ylab="Origination rates",
      xlim=4:95, plot.args=list(axes=F, cex.lab=0.8), labels.args=list(cex=0.8),
      xlab="Age (Ma)")
axis(1, cex.axis=0.8)
axis(2, cex.axis=0.8)
shades(bins$mid, as.matrix(oriDatBin[,-1]), col="darkgreen") # background
for(i in 2:ncol(oriDatBin))
  lines(bins$mid, oriDatBin[,i], lwd=0.5) # originations - BIN

# 4th panel - stages
par(mar=c(5.1, 1, 1, 4.1)) # margin adjustment
# plot
tsplot(stages, boxes="per", shading="per", ylim=c(0, 2), ylab="", xlim=4:95,
      plot.args=list(axes=F, cex.lab=0.8), labels.args=list(cex=0.8), xlab="Age (Ma)")
axis(1, cex.axis=0.8)
shades(stages$mid, as.matrix(oriDatStg[,-1]), col="darkgreen") # background
for(i in 2:ncol(oriDatStg))
  lines(stages$mid, oriDatStg[,i], lwd=0.5) # originations - stages

```



7. Conclusions

As you can see in the plots above, the estimated rates vary considerably from trial to trial, hence the variation in the final table (Table 2) summarizing the results. The different metrics have different advantages and disadvantages, although some are expected to work better than others [for instance, gap-fillers (Alroy, 2014) completely supersede the corrected three-timer rates (Alroy, 2008)]. The different subsampling procedures and stratigraphic resolutions also contribute to the variation that ultimately reflects the inherent uncertainty in our estimates.

It is apparent that some patterns are more robust than others and their uncertainty has to be taken into account when evaluating large-scale scientific hypotheses. It is highly probable that extinction and origination declined in the whole Phanerozoic, but the post-Ordovician likely featured no decline or a much shallower (lower rates in the Cenozoic are probably biologically meaningful). The most studied ‘Big Three’ (Alroy, 2008) mass extinction intervals (end-Permian, end-Triassic, end-Cretaceous) stand out clearly from the distribution of extinction rates. Equilibrial dynamics of global richness (Alroy, 2008, 2010b), on the other hand, have to be reevaluated, based on the number of different approaches we have and the increasing quantity of data.

Acknowledgments

Work on the package was funded by the Deutsche Forschungsgemeinschaft (Ko 5382/1-1, Ko 5382/1-2 and Ki 806/16-1) and is part of the Research Unit TERSANE (FO 2332). We thank all contributors of the Paleobiology Database, especially M. Clapham, A. Hendy, M. Carrano, A. Miller, M. Uhen, M. Aberhan B. Kröger, P. Wagner, M. Patzkowsky, M. Foote and J. Pálffy. We are also thankful to Na Lin, for assigning the Cambrian collections to the stages.

References

- Alroy, J. 2008. Dynamics of origination and extinction in the marine fossil record. *Proceedings of the National Academy of Science* 105:11536-11542.
- Alroy, J. 2010a. The shifting balance of diversity among major marine animal groups. *Science* 329:1191-1194.
- Alroy, J. 2010b. Geographical, environmental and intrinsic biotic controls on Phanerozoic marine diversification. *Palaeontology* 53:1211-1235.
- Alroy, J. 2014. Accurate and precise estimates of origination and extinction rates. *Paleobiology*, 40(3), 374–397.
- Alroy, J. 2015. A more precise speciation and extinction rate estimator. *Paleobiology*, 41(04), 633–639.
- Bambach, R. K., A. H. Knoll, and S. C. Wang. 2004. Origination, extinction, and mass depletions of marine diversity. *Paleobiology* 30:522-542.
- Chao, A., and L. Jost. 2012. Coverage-based rarefaction and extrapolation: standardizing samples by completeness rather than size. *Ecology* 93:2533-2547.
- Foote, M. 1994. Temporal variation in extinction risk and temporal scaling of extinction metrics. *Paleobiology*, 20(4), 424–444.
- Foote, M. 1999. Morphological diversity in the evolutionary radiation of Paleozoic and post-Paleozoic crinoids. *Paleobiology*, 25(S2), 1–115.
- Foote, M. 2000. Origination and extinction components of taxonomic diversity: General Problems. *Paleobiology* 26:74-102.
- Foote, M. 2005. Pulsed origination and extinction in the marine realm. *Paleobiology* 31:6-20.
- Good, I. J. 1953. The population frequencies of species and the estimation of population parameters. *Biometrika* 40:237-264.
- Hyndman, R. J., & Khandakar, Y. 2008. Automatic time series forecasting: the forecast package for R. *Journal of Statistical Software*, 26:1–22.
- Na, L., and W. Kiessling. 2015. Diversity partitioning during the Cambrian radiation. *Proceedings of the National Academy of Sciences* 112:4702-4706.
- Ogg, J. G., G. Ogg, and F. M. Gradstein. 2016. A concise geologic time scale: 2016. Elsevier.
- Quinn, J. F. 1983. Mass extinctions in the fossil record. *Science* 219:1239-1240.
- Raup, D. M. 1975. Taxonomic diversity estimation using rarefaction. *Paleobiology*, 1, 333–342.
- Raup, D. M., and J. J. Sepkoski. 1982. Mass extinctions in the marine fossil record. *Science* 215:1501-1503.
- Sepkoski, J. J. 1998. Rates of speciation in the fossil record. *Philosophical Transactions of the Royal Society of London B Biological Sciences*, 353:315-326.
- Sepkoski Jr, J. J. 2002. A compendium of fossil marine animal genera. *Bulletins of American Paleontology*, 363, 1-560.

Wang, X.-F. 2010. fANCOVA: Nonparametric analysis of covariance. Retrieved from <https://CRAN.R-project.org/package=fANCOVA>