

R PROFILING AND OPTIMISATION

GÜNTHER SAWITZKI

PENDING CHANGES

Warning: this is under construction.

This vignette contains experimental material which may sink down to the package implementation, or vanish.

Known issues:

- Control information may be included as special stack in raw format.
- A list of profiles may become default. Only one profiling interval value per profile.
- Nodes may be implemented as *factor*. Work-around for the R factor handling needs to be added, i.e. *factor* as a data structure.
- changing timing interval is too expensive, as *rle* is not transparent to data frames. Implement profiles as a list, with a time interval attribute per list element.

CONTENTS

Pending changes	1
Profiling facilities in R	2
LaTeX layout tools and R settings	3
1. Profiling	6
1.1. Simple regression example	7
1.1.1. R basic	8
1.1.2. Package <i>sprof</i>	12
1.2. Node attributes	18
1.2.1. Node classes	21
2. A better grip on profile information	26
2.1. The internal details	26
2.2. The free lunch	32
2.3. Cheap thrills	33
2.3.1. Trimming	36

Date: May 2013. *Revised:* Aug. 2013

Typeset, with minor revisions: October 16, 2013 from SVN *Revision* : 245 2013-10-16.

Key words and phrases. R programming, profiling, optimisation, R programming language.

*An R vignette for package *sprof*.*

URL: <http://sintro.r-forge.r-project.org/>

Private Version

2.3.2. Surgery	38
2.4. Run length	41
3. Graph package	46
4. Standard output	61
4.1. Print	61
4.2. Summary	61
4.3. Plot	62
5. More Graphs	63
5.1. Example: regression	63
5.1.1. graph package	64
5.1.2. igraph package	65
5.1.3. network package	66
5.1.4. Rgraphviz package	71
6. Template	77
Index	78
7. xxx – lost & found	81

PROFILING FACILITIES IN R

For the impatient: table 1 on page 9 and table 2 on page 11 give you the conventional information from profiling 100 runs of a simple linear regression. You get a different view of the same information by fig. 18 on page 54. The additional information we derive is summarised in table 9 on page 45 and illustrated in fig. 19 on page 55. If you want to know more, please have some patience.

R provides the basic instruments for profiling, both for time based samplers as for event based instrumentation. Information on R profiling is in section 3.2 “Profiling R code for speed” and section 3.3. “Profiling R code for memory use” of “Writing R Extensions” <http://cran.r-project.org/doc/manuals/R-expr.html>. Specific information on memory profiling is in <http://developer.r-project.org/memory-profiling.html>.

However this source of information seems to be rarely used.

Maybe the supporting tools are not adequate. The summaries provided by R reduce the information beyond necessity. Additional packages are available, but these are not sufficiently action oriented.

With package *sprof* we want to give a data representation that keeps the full profile information. Tools to answer common questions are provided. The data structure should make it easy to extend the tools as required.

The package is currently distributed at r-forge as part of the *sintro* material.

To install this package directly within R, type

```
install.packages("sprof", repos="http://r-forge.r-project.org")
```

To install the recent package from source directly within R, type

```
install.packages("sprof", repos="http://r-forge.r-project.org", type="source")
```

L^AT_EX LAYOUT TOOLS AND R SETTINGS

You may want to skip this section, unless you want to modify the vignette for your own purposes, or look at the internals.

To make sure we do not depend on packages collected along the way, we clean up packages. Calls to this function will be hidden.

Input

```
cleanpackages <- function()
{
  # make sure we do not get inherited methods or garbage here.
  try(detach("package:sna"), silent=TRUE)
  try(detach("package:igraph"), silent=TRUE)
  try(detach("package:network"), silent=TRUE)
  try(detach("package:graph"), silent=TRUE)
  try(detach("package:Rgraphviz"), silent=TRUE)
  try(detach("package:graph"), silent=TRUE)
}
```

The main library we are going to use is *sprof*.

Input

```
library(sprof)
```

We want immediate warnings, if necessary. So we set *warn* to level 1. Set *warn* to level 2 if you want to handle warnings as error.

Input

```
message("switching options(warn=1) -- immediate warning on")
options(warn=1)
```

We want a second chance on errors. So we install an error handler.

Input

```
#options(error = recover)
```

Print parameters used here:

Input

```
options(width = 72)
options(digits = 6)
```

For *str* output, we generally use these settings:

Input

```
strx <- function(x,
  max.level=2, vec.len=3,
  nchar.max=40,
  list.len=12,
  width=70, strict.width="wrap",...)
{
  cat(paste("##strx:", deparse(substitute(x)), "\n"))
}
```

```

    str(x, max.level=max.level,
        vec.len=vec.len,
        nchar.max=nchar.max,
        list.len=list.len,
        width=width, strict.width=strict.width,...)
}

```

For larger tables and data frames, we use a kludge to avoid long outputs.

```

                                Input
xcutrows <- function(df, cut, margin=5)
{
    if (!is.data.frame(df)) return(df)
    nrow <- nrow(df)
    # cut a range if it is not empty.
    # Quiet noop else.
    # Does not cut single lines.
    cutrng <- function(cutfrom,cutto, cutname="<cut>"){
        if (cutfrom<cutto){
            df[cutfrom,] <- NA
            if (!is.null(rownames(df))) rownames(df)[cutfrom] <- cutname
            if (!is.null(df$name)) df$name[cutfrom] <- ""

            cutfrom <- cutfrom+1
            df[-(cutfrom:cutto),]
        }#if
    }

    if (!missing(cut)) {
        if (length(cut) ==2)
            {df <- cutrng(cut[1],cut[2])
              return(df)}
        else {
            warning("cut ignored. Cut range must be a vector of length 2")
            return(df)}
    }

    if (length(margin)==0) return(df)
    if (length(margin)==1) margin <- c(margin, margin)
    if (length(margin)==2) {
        cut <- c(margin[1]+1,nrow-margin[2])
        df <-cutrng(cutfrom=cut[1],cutto=cut[2]);
        return(df)
    }

    #len=3: margin low, center, margin up
    delta <- (nrow-sum(margin)) %/% 2
    if (delta<2) return(df)
    c1 <- margin[1]+1
    c2 <- c1 + delta -1
    c4 <- nrow- margin[3]
    c3 <- c4 - delta
    df <- cutrng(c3,c4, cutname="<cut hi>")
    df <- cutrng(c1,c2, cutname="<cut low>")
}

```

```

    return(df)
}

```

We use the R function `xtable()` for output and L^AT_EX `longtable`. A convenient wrapper to use this in our *Sweave* source is given here. Among others, it adds `zero.print`.

ToDo: remove text
vdots from string/-
name columns.
Note: this may be
a factor. Use empty
string.

```

library(xtable)
prxt <- function (x, digits=2, cut=TRUE, caption=NULL,
                  label=NULL, zero.print=NULL, print.results=TRUE,...) {

  if (cut) {margin <- 10
    if (nrow(x)> 2*margin+3) x <- xcutrows(x, margin=margin)}

  #special sanitising for xtable
  xr <- rownames(x)
  #for (i in (1:length(xr))) {xr[i] <- sub(xr[i],"\\[","$[", fixed=TRUE)}
  #xr <- paste("$",xr,"$")
  xr <- gsub("[", "{[",xr, fixed=TRUE)
  xr <- gsub("_", "\\_",xr, fixed=TRUE)
  xr <- gsub("^", "\\^",xr, fixed=TRUE)
  rownames(x)<- xr

  pr <- print(
    xtable(x, digits=digits, caption=caption,
           label=label, ...),
    floating=FALSE,
    tabular.environment="longtable",
    caption.placement="top",
    zero.print = ".",
    NA.string="\\vdots",
    print.results=FALSE)
  # NA.string="", #NA.string="\\vdots",

  pr <- gsub( "$\\backslash$vdots","\\vdots",x=pr, fixed=TRUE)

  if(!is.null(zero.print))
    pr <- gsub( " 0 ",zero.print, x=pr, fixed=TRUE)

  if (print.results) cat(pr)
  invisible(pr)
}

```

This is to be used with `<<print=FALSE, results =tex, label=tab:prxx>>=`

The graph visualisation family is not friendly. We try to get control by using a wrapper which is at least used to the members of the *graphviz* clan. This will be used in later sections.

Input

```

plotviz <- function(x, layout="dot", main=NULL, sub=NULL,...)
{
  xid <- deparse(substitute(x))
  xsubid <- NULL
  class1 <- NULL
  if (inherits(x,"sprof")) {
    class1 <- paste0("class orig: ", paste(class(x), collapse=" "))
    xsubid <- x$info$id
    sub <- paste0(sub, " ", x$info$id)
    x <- as(adjacency(x), "graphNEL")
  }

  if (!is.null(sub)) sub <- as.character(sub)
  if (is.null(main))
    main <- paste0("plotviz( ", xid, ", ", layout, " )\n", xsubid) else
    main <- paste0(main, "\n plotviz( ", xid, ", ", layout, " )\n", xsubid)

  if (inherits(x,"Ragraph")) {
    plot(x=x, y=layout,
         cex.main=1.2,
         main=main,sub=sub,...)
  } else {
    plot(x=x, y=layout,
         attrs=list(
           node=list(cex=4, fontsize=40, shape="ellipse")),
         cex.main=1.2,
         main=main,sub=sub,...)
  }

  # title(sub = paste0(sub, " ", as.character(class(x)))
  title(sub = "xx try sub xx")
  legend("topleft",

        legend=c(class1,
                  paste0("class: ", paste(class(x), collapse=" ")),
                  paste0("layout: ",layout)),
        bg="#FFFFE040",
        seg.len=0
        )#"lightyellow" =#FFFFE0
}

```

1. PROFILING

The basic information provided by all profilers is a protocol of sampled stacks. For each recorded event, the protocol has one record, such as a line with a text string showing the sampled stack.

We use profiles to provide hints on the dynamic behaviour of programs. Most often, this is used to improve or even optimise programs. Sometimes, it is even used to understand some algorithm.

Profiles represent the program flow, which is considered to be laid out by the control structure of a program. The control structure is represented by the control graph, and this leads to the common approach to (re)construct the control graph, map the profile to this graph, and used graph based methods for further analysis. The prime example for this strategy is the GNU profiler *gprof* (see <http://sourceware.org/binutils/docs/gprof/>) which is used as master plan for many common profilers.

It is only half of the truth that the control graph can serve as a base for the profiled stacks. In R, we have some peculiarities.

lazy evaluation: Arguments to functions can be passed as promises. These are only evaluated when needed, which may be at a later time, and may then lead to insertions in the stack. So we may have information resulting from the data flow, interspersed with the control flow.

memory management: Allocation of memory, and garbage collection, may interfere and leave their traces in the stack. While allocation is closely related to the visible control flow, garbage collection is a collective effect largely out of control of the code to execute.

primitives: Internal functions may escape the usual stack conventions and execute without leaving any identifiable trace on the stack.

control structures: In R, many control structures are implemented as function. Most notably, the `apply()` family appears as function calls and can lead to cliques in the graph representation that do not correspond to relevant structures. Since these functions are well known, they can have a special treatment.

So while the stack follows an overall well known dynamics, in R there are exceptions from regularity.

The general approach, by `summaryRprof()` and others, is to reduce the profile to node information, or to consider single transitions.

We take a different approach. We take the stacks, as recorded in the profiles as our basic information unit. From this, we ask: what are the actions we need to answer our questions? Representation in graphs may come later, if they can help.

If the stacks would come from the control flow only, we could make use of the sequential nature of stacks. But since we have to live with the R specific interferences, we stay with the raw stacks.

In this presentation, we will use a small list of examples. Since *Rprof* is not implemented on all systems, and since the profiles tend to get very large, we use some prepared examples that are frozen in this vignette and not included in the distribution, but all the code to generate the examples is provided.

ToDo: rearrange
stacks? detect
order?

1.1. Simple regression example.

```
n <- 10000
x <- runif(n)
err <- rnorm(n)
```

```
y <- 2+ 3 * x + err
reg0data <- data.frame(x=x, y=y, err=err)
rm(x,y,err)
```

We will use this example to illustrate the basics. Of course the immediate questions are the variance between varying samples, and the influence of the sample size n . We keep everything fixed, so the only issue for now is the computational performance under strict iid conditions.

Still we have parameters to choose. We can determine the profiling granularity by setting the timing interval, and we can use repeated measurements to increase precision below the timing interval.

The timing interval should depend on the clock speed. Using for example 1ms amounts to some 1000 steps on a current CPU, per kernel.

If we use repeated samples, the usual rules of statistics applies. So taking 100 runs and taking the mean reduces the standard deviation by a factor 1/10.

By the usual R conventions, seconds are used as time base for parameters. However report will use ms as a time base.

Here is an example how to take a profile, using basic R. See section 1.1.2 on page 12 how to use *sampleRprof* in package *sprof* for an easier solution.

```
profinterval <- 0.001
simruns <- 100
Rprof(filename="RprofsRegressionExpl.out", interval = profinterval)
  for (i in 1:simruns) xxx<- summary(lm(y~x, data=reg0data))
Rprof(NULL)
```

We now have the profile data in a file *RprofsRegressionExpl.out*. For this vignette, we use a frozen version *RprofsRegressionExpl01.out*.

1.1.1. *R basic*. The basic R functions invite us to get a summary.

```
sumRprofRegressionExpl <- summaryRprof("RprofsRegressionExpl01.out")
#str(profile_nodes_rle, max.level=2, vec.len=3, nchar.max=40, list.len=6)
strx(sumRprofRegressionExpl)
```

```
##strx: sumRprofRegressionExpl
List of 4
 $ by.self : 'data.frame': 41 obs. of 4 variables:
  ..$ self.time : num [1:41] 0.087 0.057 0.051 0.043 0.042 0.04 0.032
    0.026 ...
  ..$ self.pct : num [1:41] 16.67 10.92 9.77 8.24 ...
  ..$ total.time: num [1:41] 0.113 0.099 0.069 0.043 0.474 0.045 0.033
    0.114 ...
  ..$ total.pct : num [1:41] 21.65 18.97 13.22 8.24 ...
 $ by.total : 'data.frame': 62 obs. of 4 variables:
```

ToDo: Can we calibrate times to CPU rate? Introduce cpu clock cycle as a time base


```

..$ total.time: num [1:62] 0.522 0.522 0.521 0.521 0.521 0.521 0.521
  0.521 ...
..$ total.pct : num [1:62] 100 100 99.8 99.8 ...
..$ self.time : num [1:62] 0.006 0 0.001 0 0 0 0 0 ...
..$ self.pct : num [1:62] 1.15 0 0.19 0 0 0 0 0 ...
$ sample.interval: num 0.001
$ sampling.time : num 0.522

```

The summary reduces the information contained in the profile to marginal statistics per node. This is provided in two data frames giving the same information, only in different order.

The file contains several spurious recordings: nodes that have been recorded only few times. It is worth noting these, but then they better be discarded. We use a time limit of 4ms, which given our sampling interval of 1ms means we require more than four observations.

```

prxt(sumRprofRegressionExpl$by.self,
      caption="summaryRprof result: by.self as final stack entry, all records",
      label="tab:prSRREbs")

```

Table 1: summaryRprof result: by.self as final stack entry, all records

	self.time	self.pct	total.time	total.pct
"lm.fit"	0.09	16.67	0.11	21.65
"{ }.data.frame"	0.06	10.92	0.10	18.97
"model.matrix.default"	0.05	9.77	0.07	13.22
"as.character"	0.04	8.24	0.04	8.24
"lm"	0.04	8.05	0.47	90.80
"summary.lm"	0.04	7.66	0.04	8.62
"structure"	0.03	6.13	0.03	6.32
"na.omit.data.frame"	0.03	4.98	0.11	21.84
"anyDuplicated.default"	0.02	4.21	0.02	4.21
"as.list.data.frame"	0.02	4.21	0.02	4.21
<cut>	:	:	:	:
"FUN"	0.00	0.19	0.01	1.34
"%in%"	0.00	0.19	0.00	0.77
"deparse"	0.00	0.19	0.00	0.38
"\$"	0.00	0.19	0.00	0.19
"as.list.default"	0.00	0.19	0.00	0.19
"as.name"	0.00	0.19	0.00	0.19
"coef"	0.00	0.19	0.00	0.19
"file"	0.00	0.19	0.00	0.19
"NCOL"	0.00	0.19	0.00	0.19
"terms.formula"	0.00	0.19	0.00	0.19

```

Input

```

ToDo: improve barplot.s. Allow vars from matrix or data frame, keep names. Use horizontal names for horizontal layout.

ToDo: fix colour. Select colours >after< sorting. Explicit override.

```
s <- sumRprofRegressionExpl$by.self$self.time
names(s) <- rownames(sumRprofRegressionExpl$by.self)
barplot_s(s, horiz=TRUE, col=rainbow(length(s)), las=1, main="self.time, by self")
```

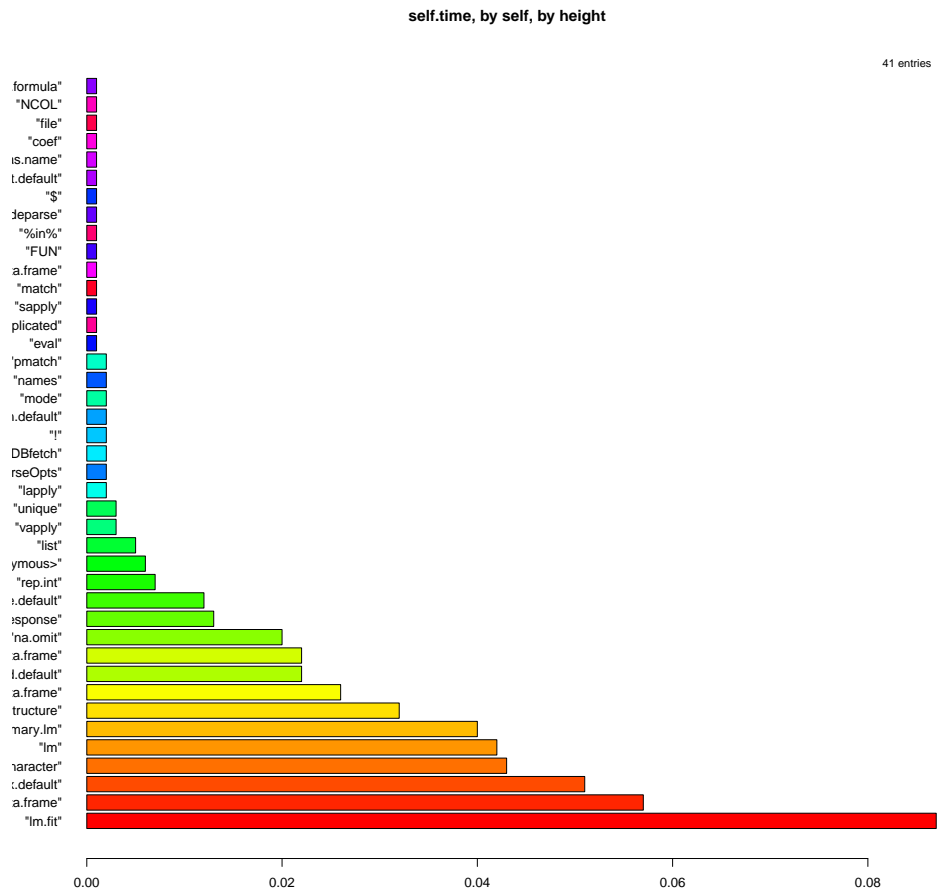


FIGURE 1. Nodes by by self.

```
prxt(
  Input
  sumRprofRegressionExpl$by.total[
    sumRprofRegressionExpl$by.total$total.time>0.004,],
  caption="summaryRprof result: by.total, total time > 4ms",
  label="tab:prSRREbt")
```

Table 2: summaryRprof result: by.total, total time > 4ms

total.time	total.pct	self.time	self.pct
------------	-----------	-----------	----------

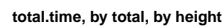
"<Anonymous>"	0.52	100.00	0.01	1.15
"Sweave"	0.52	100.00	0.00	0.00
"eval"	0.52	99.81	0.00	0.19
"doTryCatch"	0.52	99.81	0.00	0.00
"evalFunc"	0.52	99.81	0.00	0.00
"try"	0.52	99.81	0.00	0.00
"tryCatch"	0.52	99.81	0.00	0.00
"tryCatchList"	0.52	99.81	0.00	0.00
"tryCatchOne"	0.52	99.81	0.00	0.00
"withVisible"	0.52	99.81	0.00	0.00
<cut>	:	:	:	:
"as.list"	0.02	4.41	0.00	0.00
"anyDuplicated.default"	0.02	4.21	0.02	4.21
"as.list.data.frame"	0.02	4.21	0.02	4.21
"sapply"	0.01	2.68	0.00	0.19
"match"	0.01	2.11	0.00	0.19
"{ } { } .data.frame"	0.01	1.53	0.00	0.19
"{ } { }"	0.01	1.53	0.00	0.00
"rep.int"	0.01	1.34	0.01	1.34
"FUN"	0.01	1.34	0.00	0.19
"list"	0.01	0.96	0.01	0.96

Input

```

s <- sumRprofRegressionExpl$by.total$total.time
names(s) <- rownames(sumRprofRegressionExpl$by.total)
barplot_s(s, horiz=TRUE,
          col=rainbow(length(s)), las=1, main="total.time, by total")

```



1.1.2. *Package sprof.* In contrast to the common R profiling packages, in the **sprof** implementation we take a two step approach. First we read in the profile file to an internal representation. Analysis is done in later steps.

Input

Input

Output

```
$ date          : POSIXct, format: "2013-10-16 18:48:41"
$ nrnodes       : int 62
$ nrstacks      : int 50
$ nrrecords     : int 522
$ sample.interval: num 0.001
$ sampling.time  : num 0.522
$ ctllines      : chr "sample.interval=1000"
$ ctllinenr     : num 1
```

For this vignette, we change the *id* information. So in this context:

```
Input
```

```
sprof01$info$id <- "sprof01"
```

We keep this example and use the copy *sprof01* of it extensively for illustration.

```
Input
```

```
sprof01lm<-sprof01
save(sprof01lm, file="sprof01lm.RData")
rm(sprof01lm)
```

To run the vignette with a different profile, replace *sprof01* by your example. You still have the file for reference.

Package *sprof* provides a function *sampleRprof()* to take a sample and create a profile on the fly, as in

```
Input
```

```
sprof01temp <- sampleRprof(runif(10000), runs=100)
```

The basic data structure consists of four data frames. The *info* section collects global information from the input file, such as identification strings and various global matrix. The *nodes* section initially gives the same information marginal information as *summaryRprof*. The *stacks* section puts the node information into their calling context as found in the input profile file. The *profiles* section gives the temporal context. It is implemented as a list, but conceptually it is a data frame. Implementing it as a list allows run length encoding of variables, which unfortunately is not allowed by R in data frames.

ToDo: variable sampling intervals will not be supported in future versions

```
Input
```

```
strx(sprof01)
```

```
Output
```

```
##strx: sprof01
List of 4
 $ info :'data.frame': 1 obs. of 9 variables:
  ..$ id : chr "sprof01"
  ..$ date : POSIXct[1:1], format: "2013-10-16 18:48:41"
  ..$ nrnodes : int 62
  ..$ nrstacks : int 50
  ..$ nrrecords : int 522
  ..$ sample.interval: num 0.001
  ..$ sampling.time : num 0.522
```

```

..$ ctllines : chr "sample.interval=1000"
..$ ctllinenr : num 1
$ nodes : 'data.frame': 62 obs. of 7 variables:
..$ name : chr [1:62] "!" "..getNamespace" ".deparseOpts" ...
..$ self.time : num [1:62] 2 0 2 0 0 57 0 1 ...
..$ self.pct : num [1:62] 0.38 0 0.38 0 ...
..$ total.time: num [1:62] 2 1 4 26 99 99 8 8 ...
..$ total.pct : num [1:62] 0.03 0.01 0.05 0.34 1.29 1.29 0.1 0.1 ...
..$ nr_runs : num [1:62] 2 1 4 20 72 72 4 4 ...
..$ avg_time : num [1:62] 1 1 1 1.3 ...
$ stacks : 'data.frame': 50 obs. of 6 variables:
..$ nodes :List of 50
.. .. [list output truncated]
..$ refcount : num [1:50] 1 5 26 55 13 43 51 87 ...
..$ stacklength : int [1:50] 19 20 19 21 14 15 15 14 ...
..$ stackheadnodes: int [1:50] 52 52 52 52 52 52 52 52 ...
..$ stackleafnodes: int [1:50] 27 28 41 6 39 14 38 30 ...
..$ stackssrc : chr [1:50] "lazyLoadDBfetch terms model.frame.defau"|
  __truncated__ "list eval eval model.frame.default mode"|
  __truncated__ "na.omit.data.frame na.omit model.frame."|
  __truncated__ ...
$ profiles:List of 4
..$ data : int [1:522] 1 2 2 3 4 4 5 5 ...
..$ mem : NULL
..$ malloc : NULL
..$ timesRLE:List of 2
.. ..- attr(*, "class")= chr "rle"
- attr(*, "class")= chr [1:2] "sprof" "list"

```

Input

str(sprof01\$nodes)

	Output
'data.frame':	62 obs. of 7 variables:
\$ name : chr	!" "..getNamespace" ".deparseOpts" ".getXlevels" ...
\$ self.time : num	2 0 2 0 0 57 0 1 1 6 ...
\$ self.pct : num	0.38 0 0.38 0 0 ...
\$ total.time: num	2 1 4 26 99 99 8 8 4 522 ...
\$ total.pct : num	0.03 0.01 0.05 0.34 1.29 1.29 0.1 0.1 0.05 6.79 ...
\$ nr_runs : num	2 1 4 20 72 72 4 4 4 3 ...
\$ avg_time : num	1 1 1 1.3 1.38 ...

The nodes do not come in a specific order. Access via a permutation vector is preferred. This allows different views on the same data set. For example, table 4 on page 16 uses a permutation by total time, and a selection (compare to table 2 on page 11). One difference is that *sprof* uses an event count as base, usually sampled by milliseconds (ms), whereas R in general uses seconds as a base. Another difference is in two additional variables provided by *sprof*, the number of runs *nr_runs* and the average run length *avg_time*. These are discussed in section 2.4 on page 41.

Input

```

nodes <- sprof01$nodes[order(sprof01$nodes$self.time,
decreasing=TRUE),]
rownames(nodes) <- NULL
prxt(nodes[nodes$self.time>4,1:7],
digits=c(0,0,0,2,0,2,0,2),
caption="sprof result: by.self, self time > 4ms",
label="tab:prspbtself")

```

Table 3: sprof result: by.self, self time > 4ms

	name	self.time	self.pct	total.time	total.pct	nr_runs	avg_time
1	lm.fit	87	16.67	113	1.47	71	1.59
2	[.data.frame	57	10.92	99	1.29	72	1.38
3	model.matrix.default	51	9.77	69	0.90	61	1.13
4	as.character	43	8.24	43	0.56	38	1.13
5	lm	42	8.05	474	6.16	30	15.80
6	summary.lm	40	7.66	45	0.59	29	1.55
7	structure	32	6.13	33	0.43	24	1.38
8	na.omit.data.frame	26	4.98	114	1.48	63	1.81
9	anyDuplicated.default	22	4.21	22	0.29	12	1.83
10	as.list.data.frame	22	4.21	22	0.29	17	1.29
11	na.omit	20	3.83	134	1.74	71	1.89
12	model.response	13	2.49	56	0.73	42	1.33
13	model.frame.default	12	2.30	168	2.18	77	2.18
14	rep.int	7	1.34	7	0.09	7	1.00
15	<Anonymous>	6	1.15	522	6.79	3	176.00
16	list	5	0.96	5	0.07	4	1.25

At this level, it is helpful to note the expectations, and only then inspect the timing results. Since we are using a linear model, we are not surprised to see functions related to linear models on the top of the list. We may however be surprised to see functions related to data access and to character conversion very high on the list. The sizeable amount of time spent on NA handling is another aspect that is surprising.

```

Input
nodes <- sprof01$nodes[order(sprof01$nodes$total.time,
decreasing=TRUE),]
rownames(nodes) <- NULL
prxt(nodes[nodes$total.time>4,1:7],
digits=c(0,0,0,2,0,2,0,2),
caption="sprof result: by.total, total time > 4ms",
label="tab:prspbt")

```

Table 4: sprof result: by.total, total time > 4ms

	name	self.time	self.pct	total.time	total.pct	nr_runs	avg_time
1	<Anonymous>	6	1.15	522	6.79	3	176.00
2	Sweave	0	0.00	522	6.79	1	522.00

3	doTryCatch	0	0.00	521	6.78	1	521.00
4	eval	1	0.19	521	6.78	164	8.46
5	evalFunc	0	0.00	521	6.78	1	521.00
6	try	0	0.00	521	6.78	1	521.00
7	tryCatch	0	0.00	521	6.78	1	521.00
8	tryCatchList	0	0.00	521	6.78	1	521.00
9	tryCatchOne	0	0.00	521	6.78	1	521.00
10	withVisible	0	0.00	521	6.78	1	521.00
<cut>		:	:	:	:	:	:
30	vapply	3	0.57	23	0.30	16	1.44
31	anyDuplicated.default	22	4.21	22	0.29	12	1.83
32	as.list.data.frame	22	4.21	22	0.29	17	1.29
33	sapply	1	0.19	14	0.18	14	1.00
34	match	1	0.19	11	0.14	12	1.00
35	[[0	0.00	8	0.10	4	2.00
36	[[.data.frame	1	0.19	8	0.10	4	2.00
37	FUN	1	0.19	7	0.09	7	1.00
38	rep.int	7	1.34	7	0.09	7	1.00
39	list	5	0.96	5	0.07	4	1.25

ToDo: remove text
vdots from string/-
name columns

Given the sampling structure of the profiles, two aspects are common. The sampling picks up scaffold functions with a high, nearly constant frequency. And the sampling will pick up rare recordings that are near to detection range. The display functions hide these effects by default. In our example, about half of the nodes are cleared by this garbage collector.

Common rearrangements as by total time and by self time are supplied by the display functions.

`plot_nodes()`, for example, currently gives a choice of four displays for nodes, and supports trimming by default. Our profile starts with 62 nodes. The defaults cut off 34 nodes as uninformative, either because they are too rare, or ubiquitous.

Input

```
#8
oldpar <- par(mfrow=c(2,2))
plot_nodes(sprof01)
par(oldpar)
```

Basic information on node level: see fig. 3 on the next page.

Information in the time scatterplots may sometimes be more accessible when using a logarithmic scale, so this is added.

If you prefer, you can have the bar charts in horizontal layout, giving more space for labels (See “Basic information on node level - horizontal bars”: fig. 4 on page 18).

Input

```
#8
oldpar <- par(mfrow=c(2,2))
```

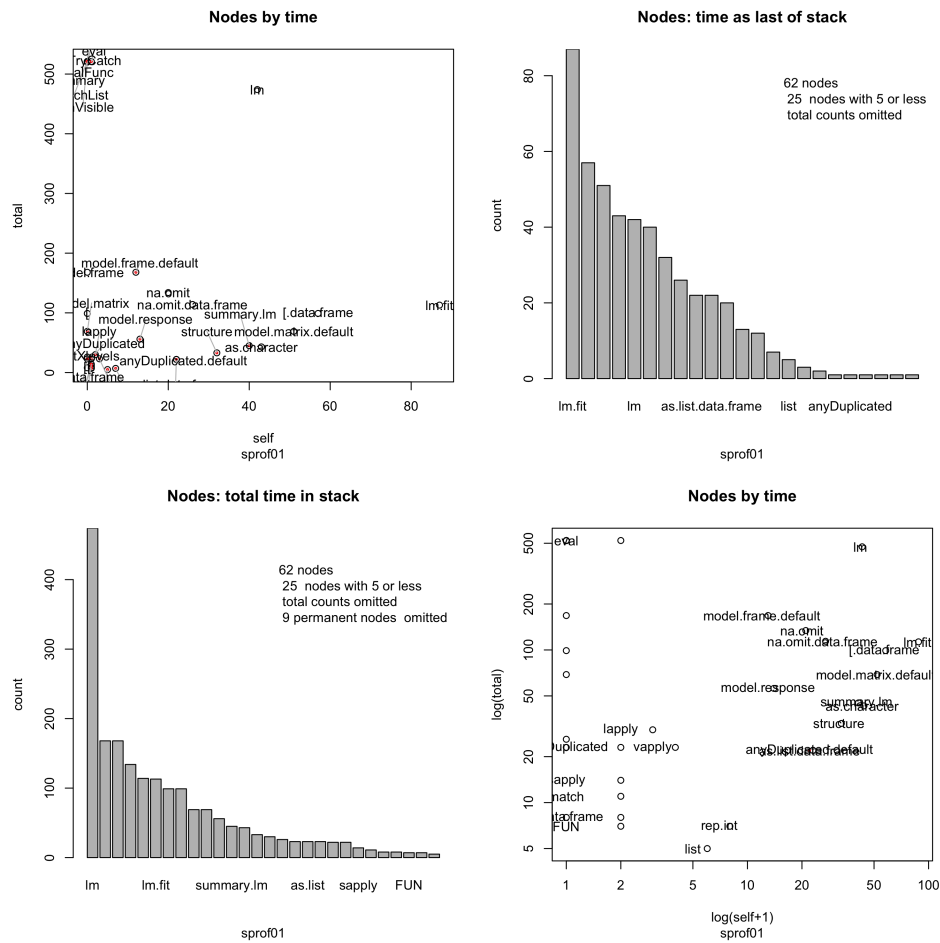
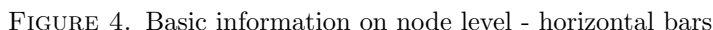



FIGURE 3. Basic information on node level

```
plot_nodes(sprof01, horiz=TRUE)
par(oldpar)
```



Keeping in mind the limits of perception, *icol* and *ilwd* should be kept small. Using 4 colours or hues seems a save choice; aboout 16 colours seem to be the maximum useful in many diplays. For line width, about 3 is a save number, and about 7 seems to be the limit.

ToDo: make call by reference function

Input

```
nodeattrs <- function(sprof,
  nodeattrs, # = c("default", "time", "runs"),
  maxcol=16, colpwr=0.5,
  maxlwd=7, lwdpow=0.5)
{
  #default
  if (missing(nodeattrs)) nodeattrs <- "default"
  icol <- rep(1,length(sprof$nodes$name))
  ilwd <- icol

  if (nodeattrs == "runs"){
    # may be skipped. if (is.null(sprof$nodes$avg_time)){
    nrl <- nodesrunlength(sprof, clean=FALSE)
    icol <- rkindex(-nrl[, "avg_time"],
      pwr=colpwr, maxindex=maxcol, ties.method="min")
    ilwd <- rkindex(nrl[, "nr_runs"],
      pwr=lwdpow, maxindex=maxlwd, ties.method="min")
  } else if (nodeattrs == "time"){
    icol <- rkindex(-sprof$nodes$self.time,
      pwr=colpwr, maxindex=maxcol, ties.method="min")
    ilwd <- rkindex(sprof$nodes$total.time,
      pwr=lwdpow, maxindex=maxlwd, ties.method="min")
  }
  data.frame(icol=icol, ilwd=ilwd)
}
```

ToDo: propagate
nodeattrs

ToDo: apply colour
to selection?

So we can add colour. To illustrate this, we encode the frequency of the nodes as colour. As a palette, we choose a heat map here.

Note: we choose a colour palette for the full data set. If we restrict a display to a selection, only a fraction of the palette may be visible. As an alternative, we could adjust the palette to focus on your selection. However the recommendation is not to spend too much energy here. Instead of modifying the display, the recommended way is to extract the core information by trimming the data set (see section 2.3.1 on page 36) and surgery (see section 2.3.2 on page 38) if necessary, and tune the displays only after these steps.

Function `plot_nodes()`, and most plotting functions of `sprof`, allow various modes to control the use of colours. The recommended mode is to add a index variable `icol` giving an index to a colour table. See `rkindex()` for a choice of modifications. By default, the colour table is the current palette. The plot functions allow to choose a different colour function plot by plot.

ToDo: add general
discussion of colours
in `sprof`

Input

```
freqrank01 <- rkindex(-sprof01$nodes$total.time,
  ties.method="random")
freqrankcol01 <- heat.colors(length(freqrank01))
```

Here is the node view, using these choices (Basic information on node level, colour by total time, see fig. 5 on the next page).

Input

```
#10
sprof01$nodes$icol <- freqranks01
oldpar <- par(mfrow=c(2,2))
plot_nodes(sprof01, col=freqrankcol01)
par(oldpar)
```

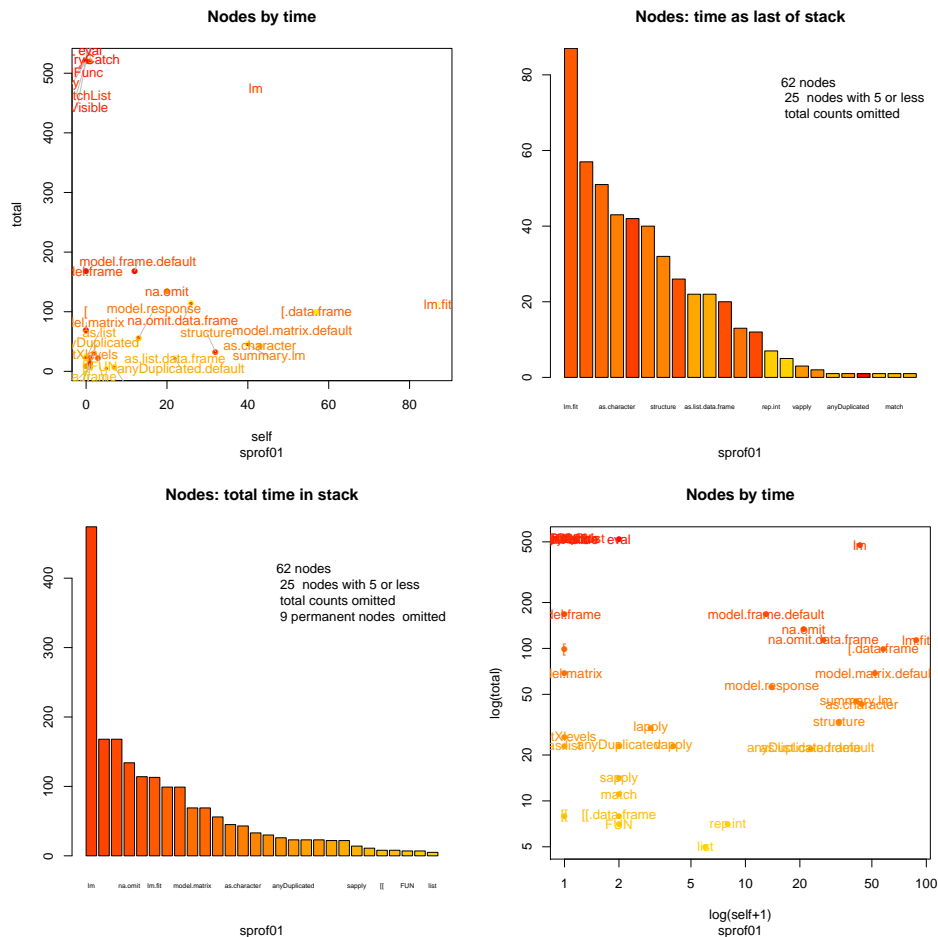


FIGURE 5. Basic information on node level, colour by total time.

Colour is considered a volatile attribute. So you may need to pay some attention to keep colour indices (and colour palettes) aligned to your context. You may want to do experiments with colour, trying to find a good solution for your visual preferences. The recommended way is to use some stable colour index (the slot `icol` is reserved for this) and use this as an index to a choice of colour palettes. So `icol` becomes a part of the data structure, and the colour palette to be used is passed as a parameter. Package *RColorBrewer* may be a helpful source for colour palettes. You find more information and references on this in the vignette “Bertin matrices” of package *bertin*, <http://bertin.r-forge.r-project.org>.

ToDo: improve
colour: support
colour in a structure

On the node level, the additional information provided by *sprof* is the information on runs. The total time is broken down to runs. Detailed information is accessible if requested. The number of runs, and the average run length are reported by default, giving an inside look into the total time reported conventionally. Making use of this information requires careful reading.

By itself, graphical inspection is not very helpful. The profiling process collapsed the real information to very few run numbers, and many ties in the average run length (fig. 6: Run length information.).

ToDo: add as additional information to classical basics

Input

```
plot(sprof01$nodes$avg_time,sprof01$nodes$nr_runs)
```

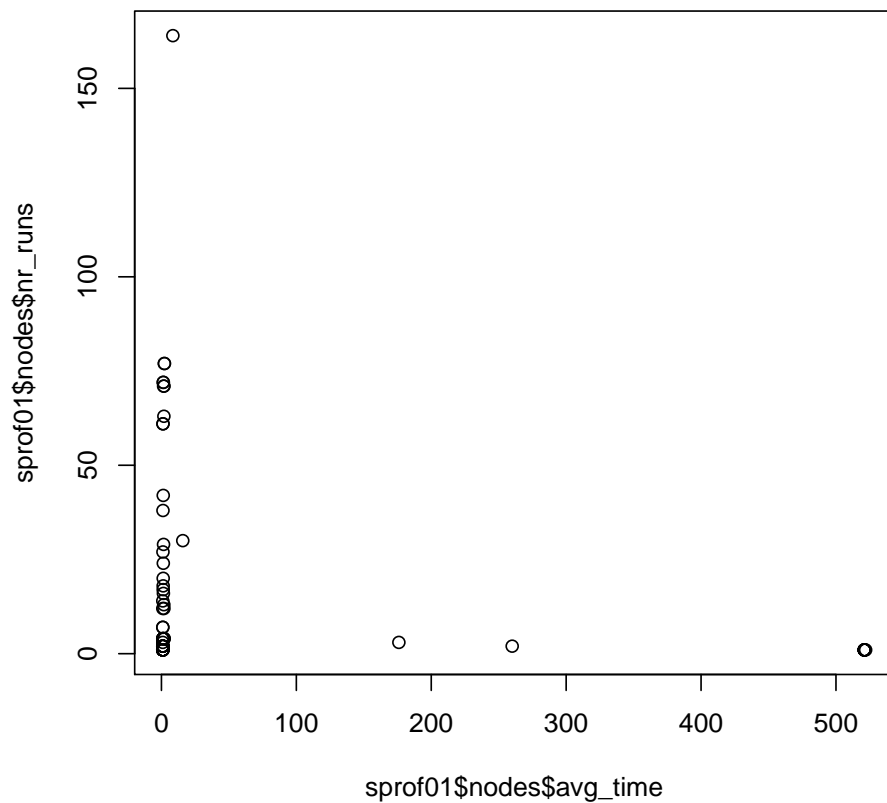


FIGURE 6. Run length information.

1.2.1. *Node classes*. We can add attributes to the plots. But we can also add attributes to the nodes, and use these in the plots. In principle, this has been always available. In *sprof*, we are making explicit use of this possibility.

The attribute `icol` is a special case which we used above. If present, it will be interpreted as an index to a colour table. For example, we can collect special well known functions in groups.

ToDo: colour by class – redo. Bundle colour index with colour?

The node information is to some part arbitrary. You may achieve the same functionality by different functions, and you will see different load in the profiles. Grouping nodes may be a mean to clarify the picture.

Grouping may also help you to focus your attention. “HOT” and “cold” may be very helpful tags. These can be used in a flexible way.

ToDo: add class by keyword

Input

```
nodekeyword0 <- function(node)
{
}
```

Input

```
nodepackages <- nodepackage(sprof01$nodes$name)
names(nodepackages) <- sprof01$nodes$name
table(nodepackages)
```

Output

```
nodepackages
<not found>      base      stats      utils
           6         41         14         1
```

Input

```
sprof01$nodes$icol <- as.factor(nodepackages)
```

Nodes by package, colours by RColorBrewer: see fig. 7 on the next page.

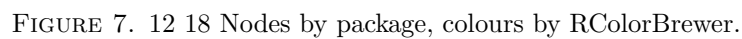
Input

```
oldpar <- par(mfrow=c(3,2))
if (require(RColorBrewer)) colpack <-
  brewer.pal(length(levels(sprof01$nodes$icol)), "Paired") else
  colpack <- rainbow(length(levels(sprof01$nodes$icol)))
plot_nodes(sprof01, which=1:6, col=colpack)
par(oldpar)
```

If you want to, you can use your own classification to group variables. We want an example which we will use later on. So we wrap it in a function.

Input

```
nodeclass <- function(sprof)
{
  # use string list to define classes
  x_apply <- c("apply", "lapply", "vapply", "sapply")
  x_as <- c("as.list", "as.data.frame", "as.list.data.frame",
           "as.character", "as.list.default", "as.name")
  # (Extend as you need it) and then use, as for example:
  ncl <- rep("x_nn", sprof$info$nrnodes)
  ncl[sprof$nodes$name %in% x_apply] <- "x_apply"
  ncl[sprof$nodes$name %in% x_as] <- "x_as"
```



```

# or use assignments on the fly
ncl[sprof$nodes$name %in%
  c("eval", "evalFunc",
    "try", "tryCatch", "tryCatchList", "tryCatchOne",
    "doTryCatch", "withVisible")
  ] <- "x_eval"
ncl[sprof$nodes$name %in%
  c("model.frame", "model.matrix.default", "model.frame.default",
    "model.response", "model.matrix", "model.response", "coef")
  ] <- "x_model"
ncl[sprof$nodes$name %in%
  c("[", "[.data.frame", "[[", "[[.data.frame") ] <- "x_acc"
ncl[sprof$nodes$name %in%
  c("lm", "lm.fit", "summary.lm")
  ] <- "x_lm"
ncl[sprof$nodes$name %in%
  c("na.omit", "na.omit.data.frame")
  ] <- "x_na"

ncl[sprof$nodes$name == "<Anonymous>"] <- "x_Anon"
ncl[sprof$nodes$name == "Sweave"] <- "x_Sweave"
ncl[sprof$nodes$name %in% c("summary", "summary.lm")] <-
  "x_summary"

return(as.factor(ncl))
}

```

```

sprof01$nodes$icol <- nodeclass(sprof01)

```

adds a sticky colour attribute. For interpretation, you should choose your preferred colour palette, for example

```

nodeclasscol <- c("red", "green", "blue", "yellow",
  "cyan", "magenta", "purple",
  "brown", "aquamarine", "pink", "violet")

```

```

# gold cyan4 aquamarine pink violet orchid hotpink salmon turquoise1

```

Nodes by user defined class, user defined colours: see fig. 8 on the facing page.

```

#8 12
oldpar <- par(mfrow=c(3,2))
plot_nodes(sprof01, which=1:6, col=nodeclasscol)
par(oldpar)

```

Nodes by user defined class: default colour selection: see fig. 9 on page 27.

```


```

ToDo: Defaults by class
ToDo: classes need separate colour palette, distinct from package or keyword.

ToDo: check/fix colour for wordcloud

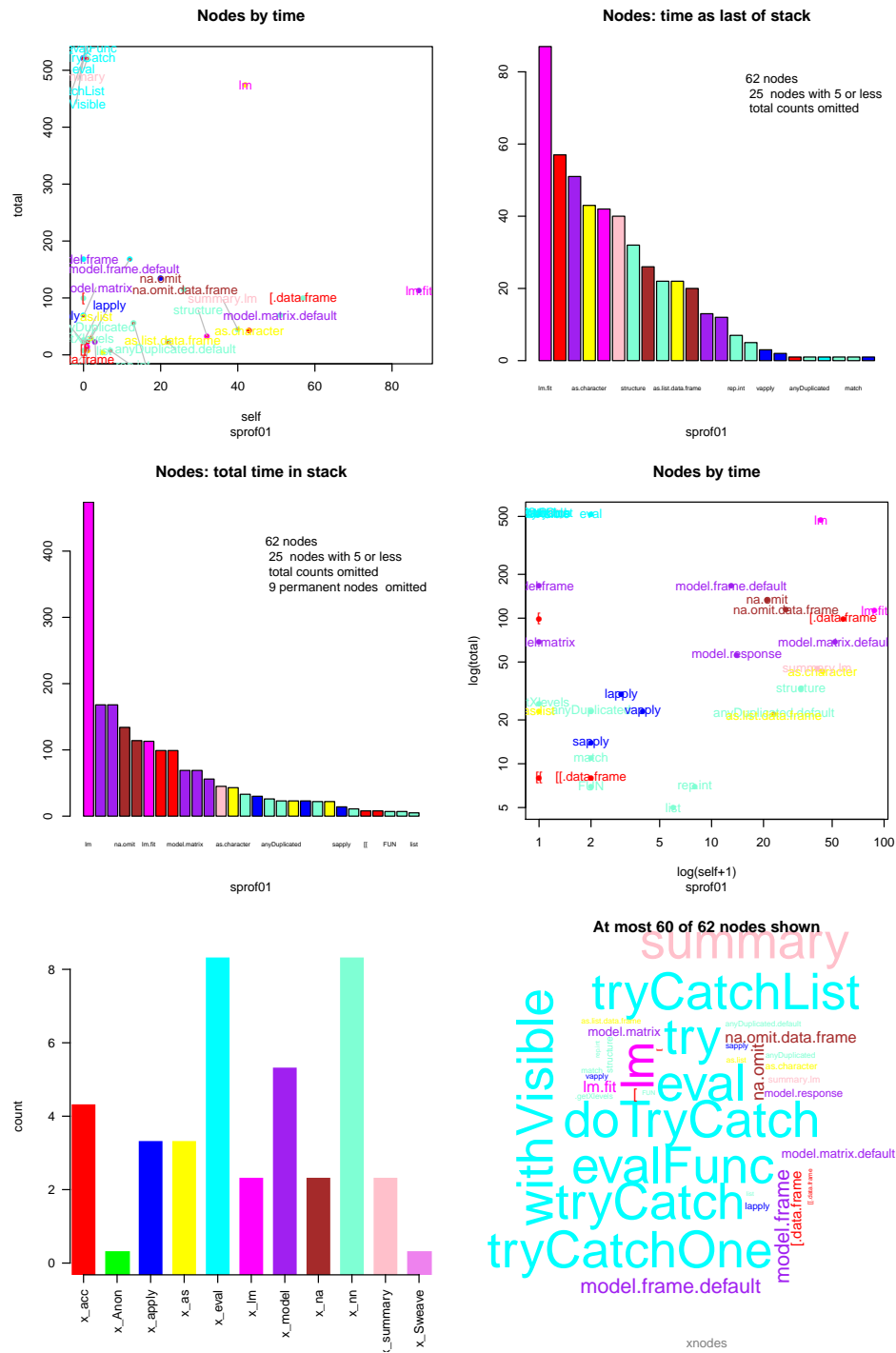


FIGURE 8. 8 12 Nodes by user defined class. User defined colours.

```
#8 12
oldpar <- par(mfrow=c(3,2))
plot_nodes(sprof01, which=1:6, las=2)
par(oldpar)
```

You can break down the frequency by classes of your choice. But beware of Simpson's paradox. The information you think you see may be strongly affected by your choices - what you see are reflections of conditional distributions. These may be very different from the global picture.

If package `wordcloud` is installed, a different view is possible. This is added in the plots above.

2. A BETTER GRIP ON PROFILE INFORMATION

The basic information provided by all profilers in R is a protocol of sampled stacks. The conventional approach is to break the information down to nodes and edges. The stacks provide more information than this. One way to access it is to use linking to pass information. This has already been used on the node level in section 1.1.2 on page 12.

ToDo: add attributes to stacks, and discuss scope

ToDo: sorting/ranging stacks

2.1. The internal details. For each recorded event, the protocol records one line with a text string showing the sampled stack (in reverse order: most recent first). The stack lines may be preceded by header lines with event specific information. The protocol may be interspersed with control information, such as information about the timing interval used.

We know that the structural information, static information as well as dynamic information, can be represented with the help of a graph. For a static analysis, the graph representation may be the first choice. For a dynamic analysis, the stack information is our first information. A stack is a connected path in the program graph. If we start with nodes and edges, we lose information which is readily available in record of stacks.

As we know that we are working with stacks, we know that they have their peculiarities. Stacks tend to grow and shrink. Subsequent events will have extensions and shrinkages of stacks (if the recording is on a fine scale), or stack sharing common stumps (if the recording is on a coarser scale). We could exploit this information, but it does not seem worth the effort.

ToDo: re-think: sort stacks

There have always been interrupts, and these show up in profiles. In R, there is a related problem: garbage collection (GC) may interfere and leave traces in the stack.

Stack information is first. The call graph is a second instance that is (re)constructed from the stack recording. The graph represents cumulated one-step information. Longer scale information contained in the stacks is lost in the graph.

Here is the way we represent the profile information:

The profile log file is sanitised:

- Control lines are extracted and recorded in a separate list.



- Head parts, if present, are extracted and recorded in a matrix that is kept line-aligned with the remainder
- Line content is standardised, for example by removing stray quotation marks etc.

After this, the sanitised lines are encoded as a vector of stacks, and references to this.

Note: after sanitising, stacks may have an empty node list. However they should not be removed as long as they may still be present in the profile.

If necessary, these steps are done by chunks to reduce memory load.

From the vector of stacks, a vector of nodes (or rather node names) is derived.

The stacks are now encoded by references to the nodes table. For convenience, we keep the (sanitised) textual representation of the stacks. (This may change.)

So far, texts are in reverse order. For each stack, we record the trailing leaf, and then we reverse order. The top of stack is now on first position.

Several statistics can be accumulated easily as a side effect.

Conceptually, the data structure consists of three tables (the implementation may differ, and is subject to change).

The profiles table is the representation of the input file. Control lines are collected in a special table. With the control lines removed, the rest is a table, one row per input line. The body of the line, the stack, is encoded as a reference to a stacks table (obligatory) and header information (optional).

The stacks table contains the collected stacks, each stack encoded as a list of references to the node table. This is obligatory. This list is kept in reverse order (root at position 1). A source line representing the stack information may be kept (optional).

The nodes table keeps the names at the nodes.

Sometimes, it is more convenient to use a simple representation, such as a matrix. Several extraction routines are provided for this, and the display routines make heavy use of this. See table 5.

ToDo: data frame conversion?

ToDo: complete matrix conversion

TABLE 5. Extraction and conversion routines

<code>profiles_matrix()</code>	incidence matrix: nodes by event
<code>stacks_matrix()</code>	incidence matrix: nodes by stack
<code>list.as.matrix()</code>	fill list to equal length and convert to matrix
<code>stackstoadj()</code>	stacks to (correspondence) adjacency matrix
<code>adjacency()</code>	sprof to (correspondence) adjacency matrix

We now can go beyond node level.

This is what we get for free from the node information on our three levels: node, stack, and profile.

ToDo: check and stabilise colour linking

See fig. 10 for a summary of nodes by stack and profile.

```
#8 rainbow
sprof01$nodes$icol <- freqrnk01; freqrnkcol <- rainbow(62)
shownodes(sprof01, col=freqrnkcol)
```

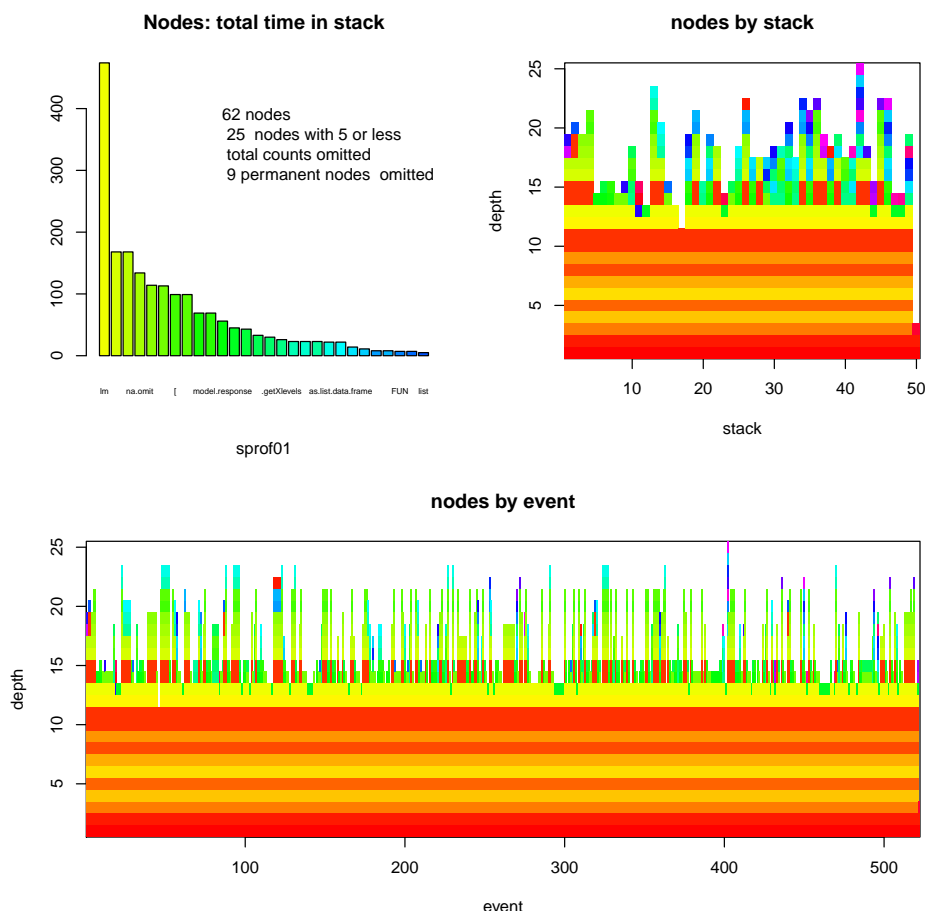


FIGURE 10. Nodes by stack and profile

The obvious message is that if seen by stack level, there are different structures. Profiling usually takes place in a framework. So at the base of the stacks, we find entries that are (almost) persistent. Then usually we have some few steps where the algorithm splits, and then we have the finer details. These can be identified using information on the stack level, but of course they are not visible on the node or edge level in a graph representation. On the stack level, we see a socket. If we want a statistic, we can look at number of different nodes by level.

Input

```
stacks_nodes <- list.as.matrix(sprof01$stacks$nodes)
nrnodes <- apply(stacks_nodes,1,function(x) {length(unique(x))})
cat("nr unique nodes per stack level\n")
```

Output

nr unique nodes per stack level

Input

nrnodes

Output

```
[1] 1 1 2 2 2 2 2 2 2 2 2 2 4 11 12 10 10 16 9 8 6 8
[23] 3 2 2
```

Nr. of unique nodes by stack level: See fig. 11.

Input

```
plot(x=nrnodes, y= 1:length(nrnodes),
      xlab="nr of unique nodes", ylab="stack level")
abline(h=2.5,col="green")
abline(h=12.5,col="green")
```

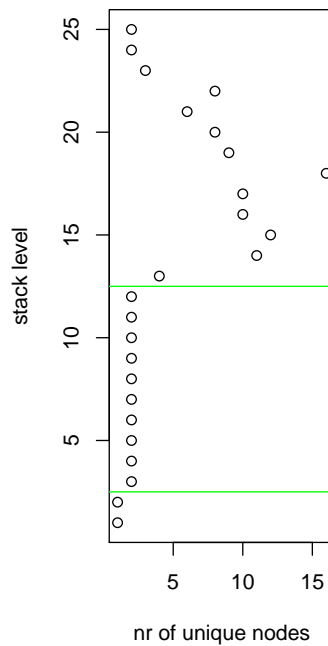


FIGURE 11. Nr of unique nodes by stack level.

ToDo: check and
synchronise

We will come to finer tools in section 2.4 on page 41 but for the moment the rough information should suffice to take a decision. In our example, it is only a matter of taste whether we cut off 12 levels, or we want to work with five components after cutting 13 levels three leaves us to start with five roots on the next level in our example.

Not so often, but a frequent phenomenon is to have some “burn in” or “fade out”. To identify this, we need to look at the profile level. The indicator to check is to whether we have very low frequency stacks at the beginning or the end of our recording. The counts to be takes as reference can be seen from the summary.

Input

```
summary(sprof01$stacks$refcount)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Output Max.
1.0	1.0	2.0	10.4	12.0	87.0

ToDo: we could do smart smoothing of the stacks here

This summary has to be taken with caution. As the program runs, the stacks are build up und teared down, and we only take random samples. So in dynamic parts, we see images with some fluctuation, as one stack may be a snapshot of an other under construction. A better information is to cut off fluctuations and use this summary as a reference.

Input

```
summary(sprof01$stacks$refcount[sprof01$stacks$refcount>2])
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Output Max.
3.00	7.75	16.00	24.30	40.50	87.00

Input

Input

```
df <- data.frame(stack=sprof01$profiles$data,
  count=sprof01$stacks$refcount[sprof01$profiles$data])
prxt(df, caption="Stacks by event: burn in/fade out",
  label="tab:margin",
  digits=c(0,0,0) )
```

Table 6: Stacks by event: burn in/fade out

	stack	count
1	1	1
2	2	5
3	2	5
4	3	26
5	4	55
6	4	55
7	5	13
8	5	13

9	6	43
10	7	51
<cut>	⋮	⋮
513	3	26
514	3	26
515	3	26
516	3	26
517	4	55
518	3	26
519	36	2
520	16	42
521	44	2
522	50	1

Here at least one recording on either side is a candidate to be off. We may have a look at the next recordings and decide to go beyond and cut off events 1 : 3 and 519 : 522.

At a closer look, we may find stack patterns (maybe marked by specific nodes) that indicate administrative intervention and rather should be handled as separators between distinct profiles rather than as part of the general dynamics. Again we may use some indicator nodes to be used as marker for special stacks. In our example, *lm* or *summary.lm* may be convenient markers.

Stable framework effects sometimes are obvious and can be detected automatically. “burn in” or “fade out” may need a closer look, and special stacks need an individual inspection on low frequency stacks. Tools for trimming are in section 2.3.1 on page 36.

2.2. The free lunch. What you have seen so far is what you get for free when using package *sprof*.

ToDo: colours. re-colour. Propagate colour to graph.

If you want to wrap up the information and look at it from a graph point of view, here is just one example. More are in section 3 on page 46 and `vrefsec:moregraph`. But before changing to the graph perspective, we recommend to see the next sections, not to skip them.

The preview, at this point, taking package *graph* as an example¹. *graph* on its side has an undocumented feature: it needs *Rgraphviz* to handle graph attributes². We have to take two steps. We extract the graph information from *sprof*. Using an adjacency matrix is a simple solution here. This is then converted to the “*graph-NEL*” format which is shared by *graph* and *Rgraphviz*. *Rgraphviz* is hidden in the use of *plot()*. So here is a bare foot approach (Call *graph* derived from profile

¹Package ‘graph’ was removed from the CRAN repository.

This package is now available from Bioconductor only.

See <http://www.bioconductor.org/packages/release/bioc/html/graph.html>.

²Package ‘Rgraphviz’ was removed from the CRAN repository.

This package is now available from Bioconductor only.

See <http://www.bioconductor.org/packages/release/bioc/html/Rgraphviz.html>.


```
#6 sprofadjNEL02
library(graph)
sprof01adjNEL <- as(adjacency(sprof01),"graphNEL")
plot(sprof01adjNEL, main="sprof01: graph layout example",
      sub=sprof01$info$id,
      attrs=list(node=list(cex=4, fontsize=40, shape="ellipse")),
      cex.main=2)
rm(sprof01adjNEL) # used for look-ahead example only
```

ToDo: log scale?

ToDo: updateRprof needs careful checking. For now, we are including long listings here to provide the necessary information

```
Input
sprof02 <- sprof01; sprof02$info$id <- "sprof02: trimmed"
```

On the stack level, we take brute force to cut off the basic stacks.

```
Input
basetrim <- 13
sprof02$stacks$nodes <- sapply(sprof02$stacks$nodes,
  function (x){if (length(x)> basetrim) x[-(1:basetrim)] })
```

We have noted burn in/fade out. This is on the profile level. Taking the big knife is not advisable, since time information and stack data must be synchronised. So we are more cautious, and instead of cutting off the stacks we replace them by an empty mark.

ToDo: should this be NA or NULL?

```
Input
sprof02$profiles$data[1:3] <- NA
sprof02$profiles$data[519:522] <- NA
```

ToDo: handle empty stacks and zero counts gracefully

ToDo: add a purge function

At this point, it is a decision whether to adapt the timing information, or keep the original information. Since this decision does affect the structural information, it is not critical. But analysis is easier if unused nodes are eliminated. The *info* section is inconsistent at this point. Another reason to call `updateRprof()`.

```
Input
strx(sprof02$info)
```

```
Output
##strx: sprof02$info
'data.frame':      1 obs. of  9 variables:
 $ id : chr "sprof02: trimmed"
 $ date : POSIXct, format: "2013-10-16 18:48:41"
 $ nrnodes : int 62
 $ nrstacks : int 50
 $ nrrecords : int 522
 $ sample.interval: num 0.001
 $ sampling.time : num 0.522
 $ ctllines : chr "sample.interval=1000"
 $ ctllinenr : num 1
```

```
Input
nodes<- sprof02$nodes[,-8]; rownames(nodes) <- NULL
prxt(nodes,
  caption="sprof02, before update",
  label="tab:sprof02info1",
  digits=c(0,0,0,2,0,2,0,2),
  zero.print=" . "
)
```

Table 7: sprof02, before update

name	self.time	self.pct	total.time	total.pct	nr_runs	avg_time
------	-----------	----------	------------	-----------	---------	----------

1	!	2	0.38	2	0.03	2	1.00
2	..getNamespace	.	0.00	1	0.01	1	1.00
3	.deparseOpts	2	0.38	4	0.05	4	1.00
4	.getXlevels	.	0.00	26	0.34	20	1.30
5	[.	0.00	99	1.29	72	1.38
6	[.data.frame	57	10.92	99	1.29	72	1.38
7	[[.	0.00	8	0.10	4	2.00
8	[[.data.frame	1	0.19	8	0.10	4	2.00
9	%in%	1	0.19	4	0.05	4	1.00
10	<Anonymous>	6	1.15	522	6.79	3	176.00
<cut>		:	:	:	:	:	:
53	terms	.	0.00	2	0.03	2	1.00
54	terms.formula	1	0.19	1	0.01	1	1.00
55	try	.	0.00	521	6.78	1	521.00
56	tryCatch	.	0.00	521	6.78	1	521.00
57	tryCatchList	.	0.00	521	6.78	1	521.00
58	tryCatchOne	.	0.00	521	6.78	1	521.00
59	unique	3	0.57	4	0.05	4	1.00
60	unlist	.	0.00	1	0.01	1	1.00
61	vapply	3	0.57	23	0.30	16	1.44
62	withVisible	.	0.00	521	6.78	1	521.00

Input

```
sprof02 <- updateRprof(sprof02)
sprof02$info$id <- "sprof02 updated"
```

```
strx(sprof02$info)
```

Output

```
##strx: sprof02$info
'data.frame':      1 obs. of  10 variables:
 $ id : chr "sprof02 updated"
 $ date : POSIXct, format: "2013-10-16 18:48:41"
 $ nrnodes : int 62
 $ nrstacks : int 50
 $ nrrecords : int 522
 $ sample.interval: num 0.001
 $ sampling.time : num 0.522
 $ ctllines : chr "sample.interval=1000"
 $ ctllinenr : num 1
 $ date_updated : POSIXct, format: "2013-10-16 18:48:47"
```

Input

```

nodes<- sprof02$nodes[,-8]; rownames(nodes) <- NULL
prxt(nodes,
  caption="sprof02, after update",
  label="tab:sprof02info2",
  #digits=c(0,0,0,2,0,2,0,2,0,2),
  digits=c(0,0,0,2,0,2,0,2),
  zero.print=" . "
)

```

Table 8: sprof02, after update

	name	self.time	self.pct	total.time	total.pct	nr_runs	avg_time
1	!	1	0.23	1	0.06	1	1.00
2	..getNamespace	.	0.00	1	0.06	1	1.00
3	.deparseOpts	2	0.46	4	0.25	4	1.00
4	.getXlevels	.	0.00	26	1.64	20	1.30
5	[.	0.00	98	6.17	71	1.38
6	[.data.frame	57	13.16	98	6.17	71	1.38
7	[[.	0.00	8	0.50	4	2.00
8	[[.data.frame	1	0.23	8	0.50	4	2.00
9	%in%	1	0.23	4	0.25	4	1.00
10	<Anonymous>	6	1.39	6	0.38	2	3.00
	<cut>	:	:	:	:	:	:
53	terms	.	0.00	1	0.06	1	1.00
54	terms.formula	1	0.23	1	0.06	1	1.00
55	try	.	0.00	.	0.00	.	0.00
56	tryCatch	.	0.00	.	0.00	.	0.00
57	tryCatchList	.	0.00	.	0.00	.	0.00
58	tryCatchOne	.	0.00	.	0.00	.	0.00
59	unique	3	0.69	4	0.25	4	1.00
60	unlist	.	0.00	1	0.06	1	1.00
61	vapply	3	0.69	23	1.45	16	1.44
62	withVisible	.	0.00	.	0.00	.	0.00

Input

Nodes by stack and profile: see fig. 13 on the facing page.

Input

```

#6 6
shownodes(sprof02)

```

2.3.1. *Trimming*. Note: trimming may be supported by the graph packages. If you are more familiar with your graph package, you may prefer to handle trimming there.

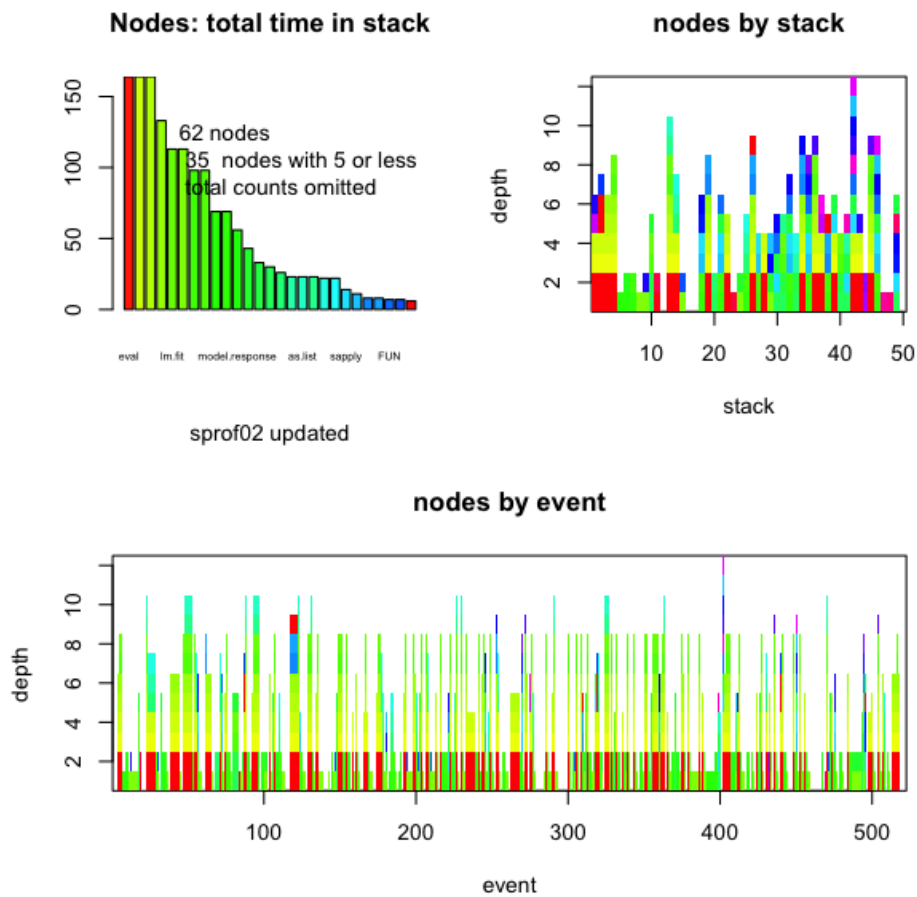


FIGURE 13. sprof02: Nodes by stack and profile

To decide about trimming, we can use the displays shown above, or we can use some statistics. Global trimming works on the stack level. To make life easy, we can imbed the stack information in a matrix and use marginals.

Input

```
stackm01 <- list.as.matrix(sprof01$stacks$nodes)
nr_uniquenodes01 <- apply(stackm01,1, function(x) { length(unique(x))})
names(nr_uniquenodes01) <- rownames(nr_uniquenodes01,FALSE,"L")
cat("Nr of unique nodes by level\n")
```

Output

```
Nr of unique nodes by level
```

Input

```
nr_uniquenodes01
```

Output																	
L1	L2	L3	L4	L5	L6	L7	L8	L9	L10	L11	L12	L13	L14	L15	L16	L17	L18
1	1	2	2	2	2	2	2	2	2	2	2	4	11	12	10	10	16
L19	L20	L21	L22	L23	L24	L25											
9	8	6	8	3	2	2											

ToDo: This section needs to be reworked

```
rm(stackm01, nr_uniquenodes01)
```

```
rsprof01Tr <- trimstacks(sprof01, level=11)
```

After trimming, it is worth to inspect the result and to see whether additional trimming is helpful. Function `roots_sprof` is handy.

Input		Output	
summary	<NA>	summary	<NA>
50	NA	50	NA

Here after trimming we have one root node. This may be aesthetically pleasing, but it is more informative to note that `summary` is a common root, and then cut it off (thus fragmenting the graph into components).

```
rsprof02Tr <- trimstacks(sprof01, level=12)
```

```
roots_sprof(sprof01, stacks=rsprof02Tr)
```

Output			
lm	lazyLoadDBfetch	summary.lm	<NA>
29	27	51	NA

Now this is revealing. We know how our test example was constructed. So we are prepared to find `lm` and `summary.lm` at prominent positions. But it is surprising to find `lazyLoadDBfetch` as dominant node, and with a frequency comparable to that of `lm`.

ToDo: move to other section

There is no statistics on profiles. Profiles are our elementary data. However we can link to our derived data to get a more informative display. For example, going one step back we can encode stacks and use these colour codes in the display of a profile.

Or going two steps back, we can encode nodes in colour, giving coloured stacks, and use these in the display of profile data.

2.3.2. *Surgery. **Note:** surgery may be supported by the graph packages. Use the implementation you are more familiar with. Surgery can also been done conveniently on the source level, using simple replace statements.*

The surgery functions will be moved to the package in the next release.

Looking at nodes gives you a point-wise horizon. Looking at edges gives you a one step horizon. The stacks give a wider horizon, typically a step size of 10 or more. The stacks we get from R have peculiarities, and we can handle with this broader perspective. These are not relevant if we look point-wise, but may become dominating if we try to get a global picture. We take a look ahead (details to come in section 3 on page 469 and have a preview how our example is represented as a graph. Left is the original graph as recovered from the edge information, right the graph after we have cut off the scaffold effects.

ToDo: cut next level

Control structures may be represented in R as function, and these may lead to concentration points. Using information from the stacks, we can avoid these by introducing substitute nodes on the stack level. For example, `lapply` is appearing in various contexts and may be confusing any graph representation. We can avoid this by replacing a short sequence.

```
"[" "lapply" ".getXlevels" -> "<.getXlevels_>"
```

Other candidates are:

```
"as.list" "vapply" "model.frame.default" -> "<model_as.list>"
```

or

```
"as.list" "vapply" "model.matrix.default" -> "<model_matrix_as.list>"
```

If the node does not exist, we want to add it to our global variable. For now, we do it using expressions on the R basic level and avoid tricks like simulating “call by reference”.

ToDo: Implement. Currently best handled on source=text level

ToDo: function addnode using “call by reference” to be added

```

                                Input
addnode <- function(nodes, newnode, warn = options("warn"))
{
  i <- match(newnode, nodes$name, nomatch=0)
  if (i==0){
    nodes$name <- as.character(nodes$name)
    nodes <- rbind(nodes,NA)
    i <- length(nodes$name)
    nodes$name[i] <- newnode
    rownames(nodes) <- nodes$name
    if (as.logical(warn))
      message("addnode: node added. An updateRprof() may be necessary.")
  }
  return(nodes)
}

```

```

                                Input
sprof03 <- sprof02; sprof03$info$id <- "sprof03: after surgery"
nodes <- addnode(sprof03$nodes, "<.getXlevels_>", warn=FALSE)
nodes <- addnode(nodes, "<model_as.list>", warn=FALSE)
nodes <- addnode(nodes, "<model_matrix_as.list>", warn=FALSE)
sprof03$nodes <- nodes

```

So far, we use factor indices only.

ToDo: xreplacenodes: improve implement

ToDo: clean up factor handling

```

xwhere <- function(oldseq, x){
  x <- unlist(x)
  firstpos <- 1; lastpos <- length(x) - length(oldseq) + 1

  if (lastpos < 1) return(0)
  l <- length(oldseq)-1
  while (firstpos <= lastpos) {
    if ( isTRUE (
      all.equal(oldseq , x[firstpos:(firstpos+1)], check.attributes=FALSE)
    ) )
      {return(firstpos) }
    firstpos <- firstpos+1
  }
  return(0)
}

```

```

Input
xreplace <- function(oldseq, newseq, x){
  #! handle multiple replacements
  wh <- xwhere(oldseq,x)
  if (wh) {
    l <- length(oldseq)-1
    if (wh>1) x1 <- c(x[1:(wh-1)], newseq) else x1 <- newseq
    if (wh+1 < length(x)) x <- c(x1, x[ -(1:(wh + 1))]) else x <- x1
    return(x)
  } else return(x)
}

```

```

Input
nodenames2index <- function(names, sprof){
  if (is.character(names))
    { sapply(names, function(x) {match(x,sprof$nodes$name)})
      #if (is.na(trimnode)) return(ts)
    } else names
}

```

```

Input
xreplacenodes <- function(sprof, oldseq, newseq)
{
  if (is.character(oldseq)) oldseq <- nodenames2index(oldseq,sprof)
  if (is.character(newseq)) newseq <- nodenames2index(newseq,sprof)
  stacks <- sprof$stacks$nodes
  sapply(stacks, function(x){xreplace(oldseq, newseq, unlist(x))})
}

```

ToDo: stacksrc,
collstacksdictrv etc.
now out of date

```

Input
sprof03$stacks$nodes <- xreplacenodes(sprof03,
  c(".getXlevels", "lapply", "["),
  "<.getXlevels_ [>")
sprof03$stacks$nodes <- xreplacenodes(sprof03,
  c("model.frame.default", "vapply", "as.list"),
  "<model_as.list>")

```



```

sprof03$stacks$nodes <- xreplacenodes(sprof03,
  c("model.matrix.default", "vapply", "as.list"),
  "<model_matrix_as.list>")

#asfactormodel(sprof03$stacks$nodes, sprof03$nodes$name)

```

Another notorious nuisance is created by shortcuts, that is aliases which maybe make programming more convenient, but may hide structure.

```

sprof03$nodes <- addnode(sprof03$nodes,
  "in_match", warn=FALSE)
sprof03$stacks$nodes <- xreplacenodes(sprof03,
  c("%in%", "match"),
  "<in_match>")

```

Now it is time to consolidate our data.

```

sprof03 <- updateRprof(sprof03)

```

The sampling process may split tight calling sequences and leave orphans.

This surgery gives no additional information on the profile data. It only helps the graph representation, by extending the one-step information given by the edges to two or more step information at some critical points.

We use a prepared sanitised version of our data set and apply our old classification.

```

sprof04 <- readRprof("RprofsRegressionExpl03.out", id="sprof04")
sprof04$nodes$icol <-nodeclass(sprof04)

```

We now get a different picture (Sanitised nodes by user defined class, user defined colours: see fig. 14 on the following page.)

```

#8 12
oldpar <- par(mfrow=c(3,2))
plot_nodes(sprof04, which=1:6, col=nodeclasscol)
par(oldpar)

```

Surgery is a means to clarify the graph structure. It should be applied with some sense. Some collusions are so intuitive that they can be ignored without surgery.

Keep in mind Simpson's paradox. If you upgrade the weights after surgery, a single node may be split to different containers, thus seemingly reducing the weights.

2.4. Run length. For a visual inspection, runs of the same node and level in the profile are easily perceived. For an analytical inspection, we have to reconstruct the runs from the data. In stacks, runs are organised hierarchically. On the root level, runs are just ordinary runs. On the next levels, runs have to be defined given (within) the previous runs. So we need `rrle()`, a recursive version of `rle`, applied

ToDo: add smoothing, remove orphans

ToDo: warn about undefined vars, e.g. `rle`, `class`, ...

ToDo: replace by `sprof03`

ToDo: fix null name

ToDo: add smart surgery with memory for attributing resources.



to the profile information. This gives a detailed information about the presence time of each node, by stack level.

Input

```
profile_nodes <- profiles_matrix(sprof03)
profile_nodes_rle <- rrle(profile_nodes, collapseNA=FALSE)
#!NA needs special case in run length handling.

strx(profile_nodes_rle, list.len=5)
```

Output

```
##strx: profile_nodes_rle
List of 11
 $ :List of 2
  ..$ lengths: int [1:365] 1 1 1 3 3 1 7 1 ...
  ..$ values : int [1:365] NA NA NA 22 39 37 30 63 ...
  ..- attr(*, "class")= chr "rle"
 $ :List of 2
  ..$ lengths: int [1:413] 1 1 1 3 1 1 1 1 ...
  ..$ values : int [1:413] NA NA NA 22 NA NA 14 38 ...
  ..- attr(*, "class")= chr "rle"
 $ :List of 2
  ..$ lengths: int [1:431] 1 1 1 3 1 1 1 1 ...
  ..$ values : int [1:431] NA NA NA 35 NA NA NA NA ...
  ..- attr(*, "class")= chr "rle"
 $ :List of 2
  ..$ lengths: int [1:441] 1 1 1 3 1 1 1 1 ...
  ..$ values : int [1:441] NA NA NA 36 NA NA NA NA ...
  ..- attr(*, "class")= chr "rle"
 $ :List of 2
  ..$ lengths: int [1:456] 1 1 1 3 1 1 1 1 ...
  ..$ values : int [1:456] NA NA NA 40 NA NA NA NA ...
  ..- attr(*, "class")= chr "rle"
[list output truncated]
```

Input

On a given stack level, the run length is the best information on the time used per call, and the run count of a node is the best information on the number of calls. So this is a prime starting point for in-depth analysis.

If you need it, you can represent the run length information by level as a matrix. This is expanding a sparse matrix to full and should be avoided.

Input

```
profile_nodes_rlearray <- nodesprofile(sprof03)
strx(profile_nodes_rlearray)
```

Output

```
##strx: profile_nodes_rlearray
num [1:66, 1:11, 1:7] 0 1 0 3 0 0 0 0 ...
- attr(*, "dimnames")=List of 3
```

ToDo: keep as factor. This is a sparse cube with margins node, stack level, run length. Nodes are mostly concentrated on few levels.

ToDo: Warning: data structure still under discussion

ToDo: hack. replace by decent vector/array based implementation

```

..$ node : chr [1:66] "!" "..getNamespace" ".deparseOpts" ...
..$ level : chr [1:11] "1" "2" "3" ...
..$ run_length: chr [1:7] "1" "2" "3" ...

```

ToDo: add summary for NA
ToDo: add marginals and conditionals.
 Provide function node_summary.

This allows us to extract marginal from `provlev[node, level, run length]`.

Input

```

nn <- profile_nodes_rlearray["model.frame", , ]
print.table(addmargins(nn), zero.print = ".")

```

Output

```

      run_length
level  1  2  3  4  5  6  7 Sum
  1    .  .  .  .  .  .  .
  2    .  .  .  .  .  .  .
  3   40 17  7  4  2  6  1  77
  4    .  .  .  .  .  .  .
  5    .  .  .  .  .  .  .
  6    .  .  .  .  .  .  .
  7    .  .  .  .  .  .  .
  8    .  .  .  .  .  .  .
  9    .  .  .  .  .  .  .
 10    .  .  .  .  .  .  .
 11    .  .  .  .  .  .  .
Sum  40 17  7  4  2  6  1  77

```

Input

```

amt <- nodesrunlength(sprof03)
prxt(amt,
     caption=paste0("Marginal statistics on nodes by run length, ",
                    "sorted by total time used"),
     label="tab:pramt4",
     digits=c(rep(0,dim(amt)[2]) ,2),
     zero.print=" . ") #dim(amt)[2]-1, +1 for rownames

```

Table 9: Marginal statistics on nodes by run length, sorted by total time used

	1	2	3	4	5	6	7	nr_runs	total_time	avg_time
eval	86	34	14	8	4	12	2	160	334	2.09
model.frame	40	17	7	4	2	6	1	77	164	2.13
model.frame.default	42	15	7	3	2	5	1	75	152	2.03
na.omit	46	10	5	4	1	4	1	71	133	1.87
lm.fit	46	18	3	1	1	1	1	71	113	1.59
na.omit.data.frame	43	7	4	5	.	4	.	63	113	1.79
{[]}.data.frame	59	4	3	4	.	1	.	71	98	1.38
{[]}	45	4	3	3	.	1	.	56	80	1.43
model.matrix	55	4	2	61	69	1.13
model.matrix.default	52	3	1	56	61	1.09
model.response	35	3	3	.	.	1	.	42	56	1.33
as.character	34	3	1	38	43	1.13
structure	21	1	.	1	.	1	.	24	33	1.38

anyDuplicated	10	.	.	2	1	.	.	13	23	1.77
anyDuplicated.default	9	.	.	2	1	.	.	12	22	1.83
as.list.data.frame	15	1	.	.	1	.	.	17	22	1.29
<.getXlevels_{>	14	.	.	1	.	.	.	15	18	1.20
sapply	14	14	14	1.00
lapply	12	12	12	1.00
<model_as.list>	7	.	.	.	1	.	.	8	12	1.50
match	9	9	9	1.00
.getXlevels	3	1	1	5	8	1.60
{ } { }	3	.	.	.	1	.	.	4	8	2.00
{ } { }.data.frame	3	.	.	.	1	.	.	4	8	2.00
<model_matrix_as.list>	6	1	7	8	1.14
FUN	7	7	7	1.00
rep.int	7	7	7	1.00
<Anonymous>	1	.	.	.	1	.	.	2	6	3.00
.deparseOpts	4	4	4	1.00
simplify2array	4	4	4	1.00
unique	4	4	4	1.00
as.list	3	3	3	1.00
list	3	3	3	1.00
vapply	.	.	1	1	3	3.00
<in_match>	3	3	3	1.00
deparse	2	2	2	1.00
mode	2	2	2	1.00
names	2	2	2	1.00
pmatch	2	2	2	1.00
!	1	1	1	1.00
..getNamespace	1	1	1	1.00
%in%	1	1	1	1.00
\$	1	1	1	1.00
as.list.default	1	1	1	1.00
as.name	1	1	1	1.00
coef	1	1	1	1.00
lazyLoadDBfetch	1	1	1	1.00
mean	1	1	1	1.00
mean.default	1	1	1	1.00
NCOL	1	1	1	1.00
paste	1	1	1	1.00
terms	1	1	1	1.00
terms.formula	1	1	1	1.00
unlist	1	1	1	1.00

See table 9: Marginal statistics on nodes by run length.

From the summary information, **nr_runs** and **avg_time** are included in the node information by default. We can use this information to enhance graphical displays. See: nodes marked by run length and run count, fig. 19 on page 55.

ToDo: hack. keep length in nodesrun-length

`eval` is the base of all evaluation in R, so we should not be surprised to find it at the top of the list. The next two entries, `model.frame` and `model.frame.default` are seen 77 times. We know that in our example we had a loop with 100 repetitions, so a frequency of 100 would not be surprising. A closer look shows that both functions occur as isolated events 40 times, but frequency decreases with run length until we see another high at run length 6. We are sampling, and of course sampling can miss some function calls. But this pattern is the opposite. This is what if occurs if we do not miss a function call, but we miss a gap. So several function calls are joined and appear as some longer runs, typically a multiple of the original run length.

We can improve the hit rate by increasing the sampling rate. But of course this is at the expense of using more space for the log files, and increased time for the overhead. In our case, 1 ms seems to be a good compromise, but 0.1 ms might be another feasible choice.

As we walk down the list, `na.omit` is next, followed closely by `na.omit.data.frame`. In our problem, there are no missing data. But R does not have something like a “vanilla mode”. The overhead used for the handling of potentially missing values would need a restructuring of the linear model algorithm, e.g. by introducing a “has.na” flag.

Walking further down, we see other candidates such as `as.character` that may be avoidable in this problem.

ToDo: table: node
#runs min median
run length max

3. GRAPH PACKAGE

What we have achieved so far can be seen from the graph representations. For our purposes, an edge table is most convenient. To allow for edge attributes, we can use an R `data.frame` as provided by

We can make use of any graph mapping package. Unfortunately, each seem to have its own calling convention. So we have to do some translation.

Input

```
#12 6
library(graph)
oldpar <- par(mfrow=c(1,2))
  plotviz(sprof01)
  plotviz(sprof02)
par(oldpar)
```

See fig. 15 on the next page for a comparison before and after trimming. The scaffold effect are removed from the picture on the right side. This cuts off the uninformative spine, and induces minor changes in the body of the graph. You can do additional trimming, if you want to.

Input

```
#12 6
library(graph)
oldpar <- par(mfrow=c(1,2))
  plotviz(sprof03)
```

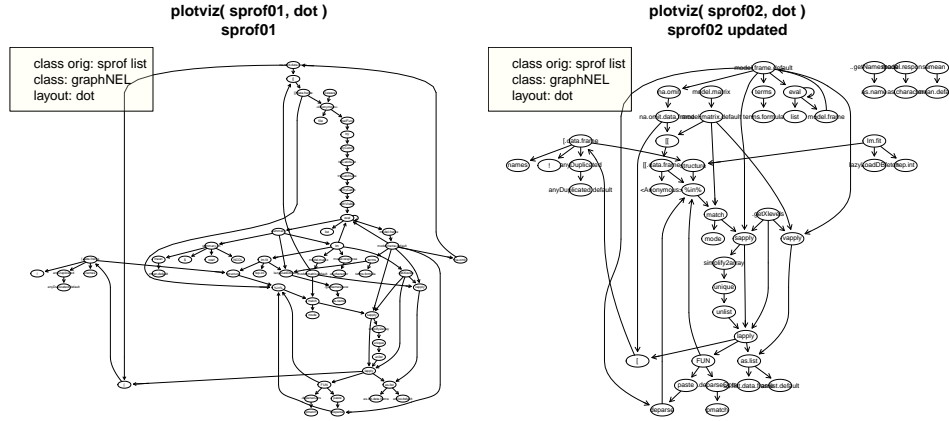


FIGURE 15. Sprof graph, before and after trimming.

```
plotviz(sprof04)
par(oldpar)
```

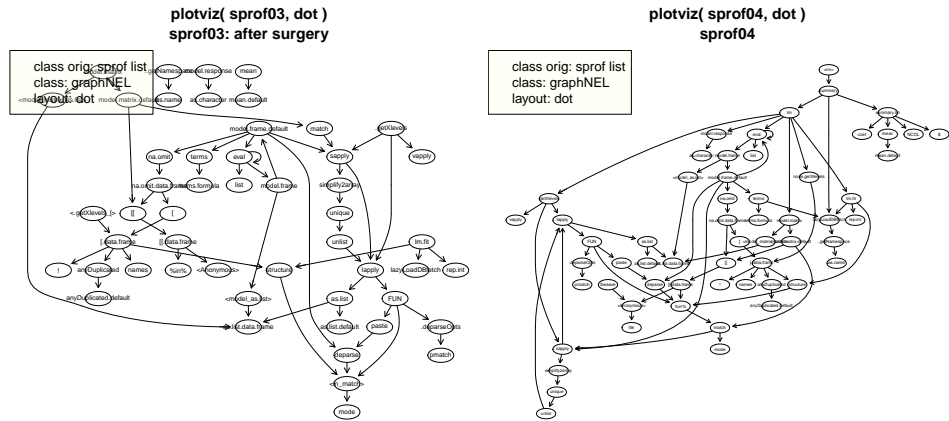


FIGURE 16. Sprof graph, two variants of surgery.

R is function based, and control structures in general are implemented as functions. In a graph representation, they appear as nodes, concentrating and seeding to unrelated paths. We can detect these on the stack level and replace them by surrogates, introducing new nodes. This is a case for surgery. See fig. 16 for two variants of surgery. While it is possible to apply surgery on the editor level, the experience is that it is worth to encode the surgery rules. So *Sprof03* will be preferred over *Sprof04*.

Some annoying concentrations remain in the graph, in particular nodes from the `apply` family. But these are so familiar that it seems to be easy to skip them visually. We did not take the pain to resolve them.

ToDo: cut top levels

Now the structure becomes obvious. Cutting off may be taken two levels deeper. This would completely separate the `lm()` branch from the `summary.lm()`. In the `lm()` branch, there are tree nodes (`lapply()`, `%in%` and `[.data.frame]`) that are cases for additional surgery to avoid focusing affects in the graph display. It is your choice to remove them, or live with them.

We know how to create standard graph displays from this. The next step is to encode additional information we have from the profiles as attribute to the graph.

The derived edge frequency is the first bit of information. Implicitly, it can be used as weight in the graph placement routines. We make this explicit by giving a choice whether to use it or not. Irrespective of this choice, we encode reference counts as line width of the edges.

The functions in this chapter are not included in the package to avoid dependency of `sprof` on `graph` and other graph packages.

Input

```
library(graph)
library(Rgraphviz)
```

This is a common routine for the `graph` and `Rgraphviz` package.

Input

```
as_graphNEL_sprof <- function(sprof, weight=TRUE){

  a04<-adjacency(sprof)

  rnames <- rownames(a04)

  if (!weight) {
    dimold <- dim(a04); a04 <- as.numeric(a04); dim(a04) <- dimold
    rownames(a04)<- rnames; colnames(a04)<- rnames;
  } #! define lwd first

  el04 <- edgedf(a04)
  el04$lwd <- rindex(el04$count, maxindex=7, ties.method="min")

  a04NEL <- as(a04,"graphNEL")
  nodeDataDefaults(a04NEL, "shape") <- "ellipse"
  nodeDataDefaults(a04NEL, "cex") <- 0.6
  nodeDataDefaults(a04NEL, "weight") <- 1
  nodeDataDefaults(a04NEL, "fill") <- "green"
  nodeDataDefaults(a04NEL, "col") <- "yellow"

  a04NEL <- layoutGraph(a04NEL)

  nodeRenderInfo(a04NEL) <- list(shape="ellipse")
  nodeRenderInfo(a04NEL) <- list(cex=0.6, shape="ellipse")
}
```



```

nodeRenderInfo(a04NEL) <- list(weight=1)
#nodeRenderInfo(a04NEL) <- list(color="yellow")
nodeRenderInfo(a04NEL) <- list(fill="yellow", col="blue")

edgeDataDefaults(a04NEL,"lwd") <- 1
edgeDataDefaults(a04NEL,"col") <- "grey"

#nodeRenderInfo(a04NEL) <- list(weight=1)

#edgeRenderInfo(a04NEL) <- list(lwd=e104$lwd)
#edgeRenderInfo(a04NEL)$lwd <- e104$lwd
for (i in 1:length(e104$lwd))
{edgeRenderInfo(a04NEL)$lwd[i] <- e104$lwd[i]}
a04NEL
}

```

As `as_graphNEL_sprof`, the following function is not included in the package to avoid dependency of `sprof` on `graph` and other graph packages.

```

# source('~/projects/rforge/sintro/pkg/sprof/R/pgn.R', chdir = TRUE)
plot_graphNEL_sprof <- function(sprof,
  layoutType = "dot",
  nodeattrs, # = c("default", "time", "runs"),
  fill_list = NULL,
  lwd_list = NULL,
  srclegend="topleft",
  attrlegend="bottomright",
  main = NULL, sub = NULL,...)
{
  main1 <- deparse(substitute(sprof))
  xsubid <- NULL
  y<-NULL
  class1 <- NULL

  if (missing(nodeattrs)) nodeattrs <- "default"
  #nodeattrs <- match.arg(nodeattrs)
  # cat(nodeattrs)

  if (inherits(sprof, "sprof")) {
    class1 <- paste0("class orig: ", paste(class(sprof), collapse=" "))
    main1 <- sprof$info$id
    graphNEL <-
      as_graphNEL_sprof(sprof, weight=FALSE)
    y <- layoutType
  } else graphNEL <- sprof

  if (!inherits(graphNEL, "graphNEL")) warn("graphNEL data structure needed")

  if (!is.null(sub)) sub <- as.character(sub)

  main=paste0(main, "\n plot_graphNEL_sprof( ",
    main1 , " , ",

```

```

        layoutType, " )\n", xsubid)

# functions
legattributes <- function(legnd=NULL,
                          nodecol=NULL, node1wd=NULL,
                          edgcol=NULL, edglwd="frequency")
{
  #
  llegnd<- legnd
  if (!is.null(nodecol)) {
    nodecol <- paste0("node col: ",nodecol, collapse="")
    if (is.null(llegnd)) llegnd <- nodecol else
      llegnd <- paste0(llegnd,"\n", nodecol,collapse="")
  }
  if (!is.null(node1wd)) {
    node1wd <- paste0("node lwd: ", node1wd, collapse="")
    if (is.null(llegnd)) llegnd <- node1wd else
      llegnd <- paste0(llegnd,"\n", node1wd, collapse="")
  }
  if (!is.null(edgcol)) {
    edgcol <- paste0("edge col: ", edgcol, collapse="\n")
    if (is.null(llegnd)) llegnd <- edgcol else
      llegnd <- paste0(llegnd,"\n", edgcol,collapse="")
  }
  if (!is.null(edglwd)) {
    edglwd <- paste0("edge lwd: ", edglwd, collapse="")
    if (is.null(llegnd)) llegnd <- edglwd else
      llegnd <- paste0(llegnd, "\n",edglwd,collapse="")
  }
  if (!is.null(attrlegend)){
    if (!is.null(llegnd))
    {
      llegnd=paste0(llegnd,"\n")
    }

    legend(attrlegend,
            legend= llegnd,
            bg="#FFFFE0",
            seg.len=0,
            bty="n",
            text.font=3)
    #bg="#0000"
  }
} # legattributes

graphNEL <- layoutGraph(graphNEL, layoutType=layoutType)

nodeDataDefaults(graphNEL, "shape") <- "ellipse"

if (nodeattrs == "runs"){
  amt <- nodesrunlength(sprof, clean=FALSE)

  #sprof$nodes$self.time -> fill
  fill_list <- heat.colors(12)[

```

```

        rkindex(-amt[, "avg_time"],
                pwr=0.5, maxindex=12, ties.method="min")]
names(fill_list) <- sprof$nodes$name

#sprof$nodes$total.time -> lwd
lwd_list <- rkindex(amt[, "nr_runs"],
                  pwr=0.5, maxindex=7, ties.method="min")
names(lwd_list) <- sprof$nodes$name

#strx(nodeRenderInfo(graphNEL))
nDD3 <- (nodeRenderInfo(graphNEL))

nodeDataDefaults(graphNEL, "shape") <- "ellipse"

nodeRenderInfo(graphNEL) <- list(lwd=lwd_list,
                                fill=fill_list,
                                col="#0000FF80",
                                shape="ellipse",
                                weight=1)
#strx(nodeRenderInfo(graphNEL))
} else if (nodeattrs == "time"){
  # node attributes
  #sprof$nodes$self.time -> fill
  fill_list <- heat.colors(12)[
    rkindex(-sprof$nodes$self.time,
            pwr=0.5, maxindex=12, ties.method="min")]
names(fill_list) <- sprof$nodes$name

  #sprof$nodes$total.time -> lwd
  lwd_list <- rkindex(sprof$nodes$total.time,
                    pwr=0.5, maxindex=7, ties.method="min")
names(lwd_list) <- sprof$nodes$name

  #strx(nodeRenderInfo(graphNEL))
  nDD3 <- (nodeRenderInfo(graphNEL))
} else {
  # graphNEL <- layoutGraph(graphNEL, layoutType=layoutType)
  fill_list <- NULL
  lwd_list <- NULL
}
nodeDataDefaults(graphNEL, "shape") <- "ellipse"

nodeRenderInfo(graphNEL) <- list(
  lwd=lwd_list,
  fill=fill_list,
  col="#0000FF80",
  shape="ellipse",
  weight=1)

# edge attributes
# strx(edgeRenderInfo(graphNEL))
nER01 <- edgeRenderInfo(graphNEL)
edgeRenderInfo(graphNEL) <- list(col= "#80808080")

```

```

# strx(edgeRenderInfo(graphNEL))

renderGraph(graphNEL)

if (!is.null(srclegend))
  legend(srclegend,
        legend=c( class1,
                   paste0("class: ",
                           paste(class(graphNEL), collapse=" ")),
                   paste0("layout: ", layoutType)),
        bg="#FFFFE040",
        seg.len=0
        )#"lightyellow" =#FFFFE0

  if (nodeattrs == "runs"){
    title(main=paste0("nodes (run length)\n", main)          )
    legattributes(nodecol="avg run len, pwr=0.5",
                  nodelwd="nr runs, pwr=0.5")
  } else if (nodeattrs == "time"){
    title(main=paste0("nodes (time)\n", main))
    legattributes( nodecol="self.time, pwr=0.5",
                  nodelwd="total.time, pwr=0.5")
  } else {
    title(main=main)
    legattributes( )
  }
  title( sub=spref$info$id, col.sub=grey(0.5))
  invisible(graphNEL)
}# %: plot_graphNEL_sprof(sprof03)

_____ Input _____
plot_graphNEL_sprof(sprof03)

```

ToDo: remove
global colour; imple-
ment local colour

ToDo: merge with
as_graphNEL_sprof

Rgraphviz/graph basic plot: see fig. 17 on the next page.

To use attributes on nodes and edges, we need *Rgraphviz*.

Rgraphviz/graph plot with attributes: see fig. 18 on page 54. This plot gives us the traditional view, highlighting nodes and edges with overall time presence.

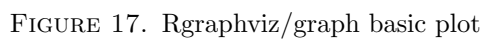
```

_____ Input _____
#6 6
plot_graphNEL_sprof(sprof03, nodeattrs="time")

```

We should not overload the plot. We could use colour encoding for the edges, but this would conflict visually with the colour encoding of the nodes. We could use different shapes for classes of nodes, but then we would need an additional display to explain the shape information.

But within the choice of attributes used, we still can select the information shown. To close this round, instead of showing the node information from the rough summary, we can show the information from the run length discussed in section 2.4 on page 41.



#6 6
Input

```
library(Rgraphviz)
plot_graphNEL_sprof(sprof03, nodeattrs="runs")
```

The bottom line is: graph layout is best left to graph representation packages. You can help by trimming and surgery. As far as attributes are concerned, the current

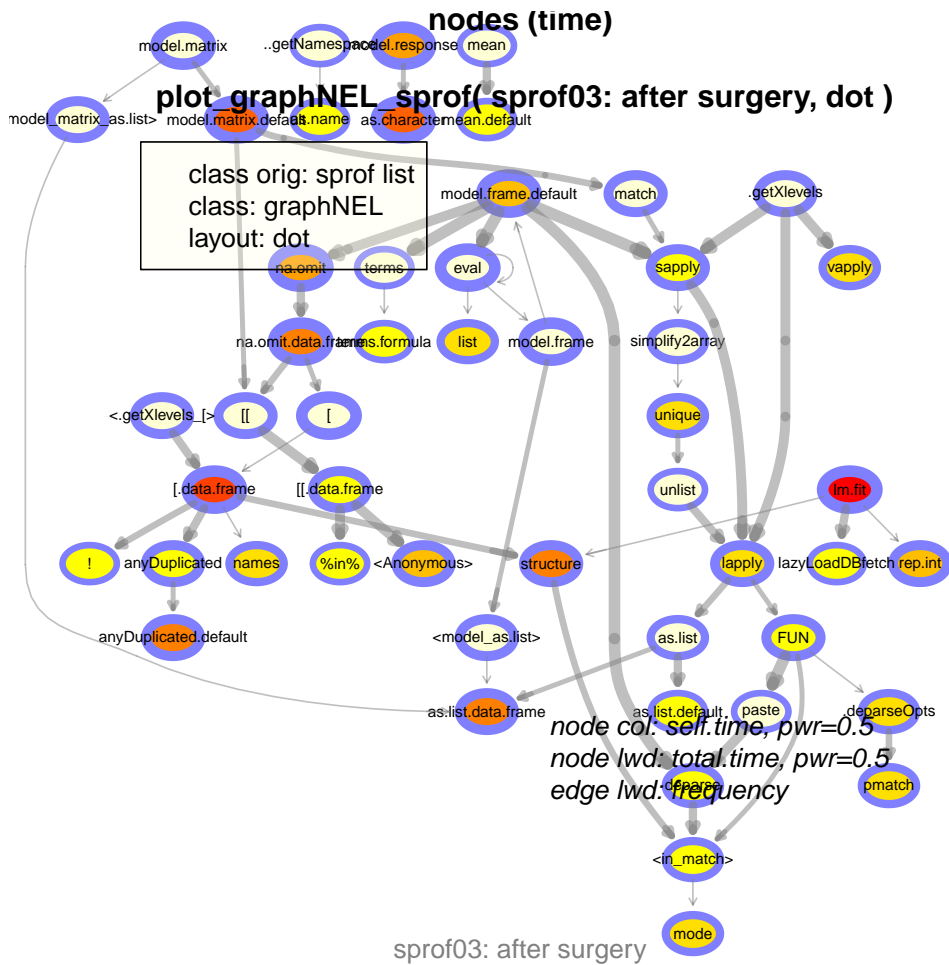


FIGURE 18. Rgraphviz/graph plot with attributes

recommendation is to start with run time lengths and run time frequencies by node and level:

```
plot_graphNEL_sprof(x, nodeattrs = "runs")
```

Here for comparison is a run-down of the previous stages of our example, using this plot.

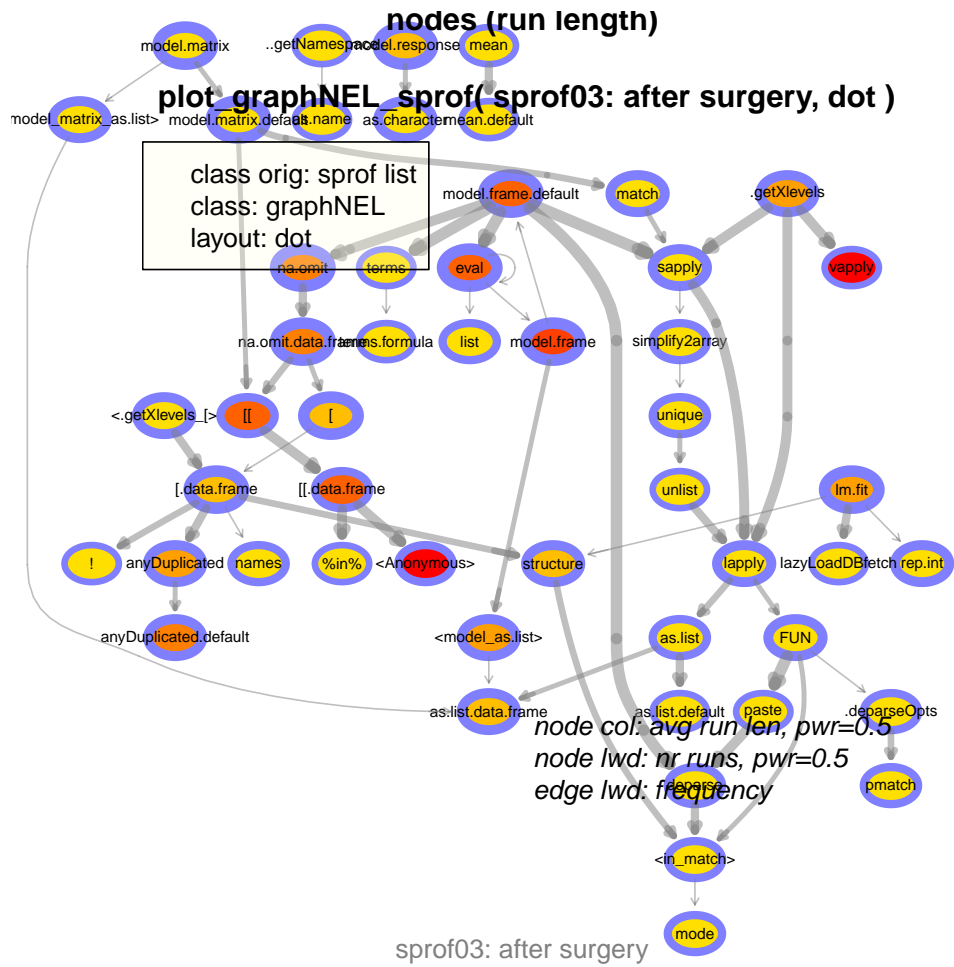
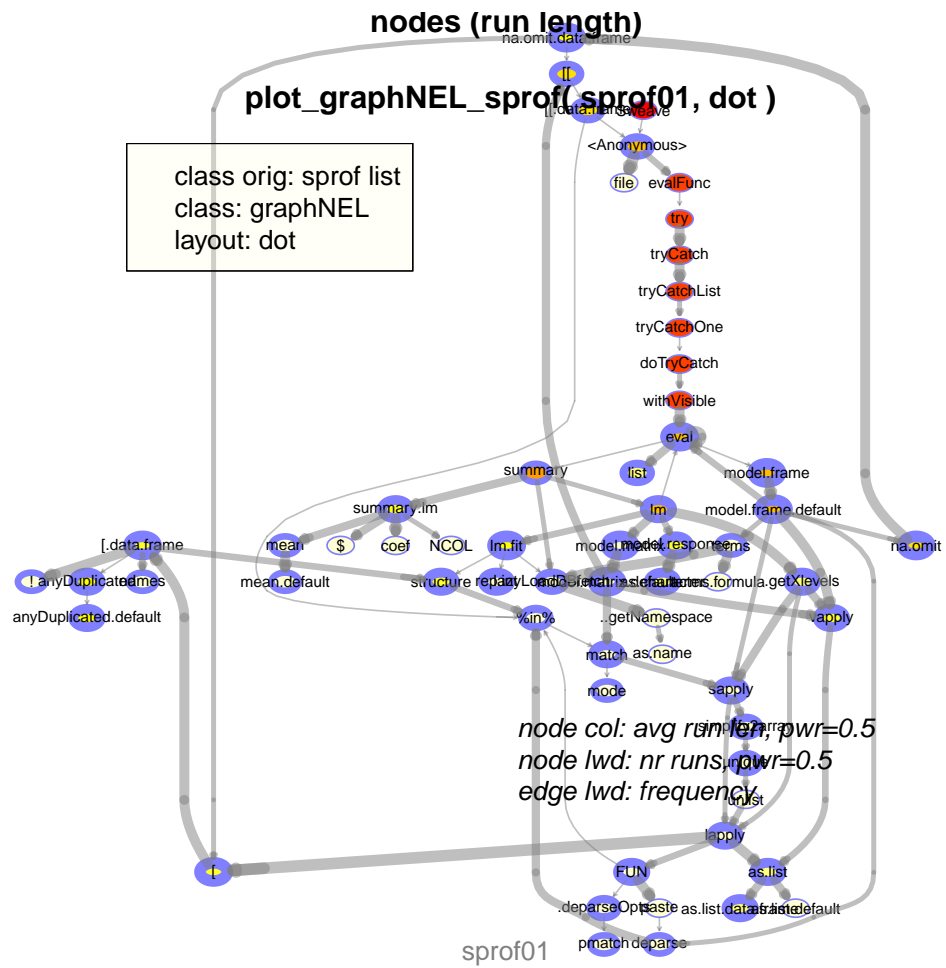


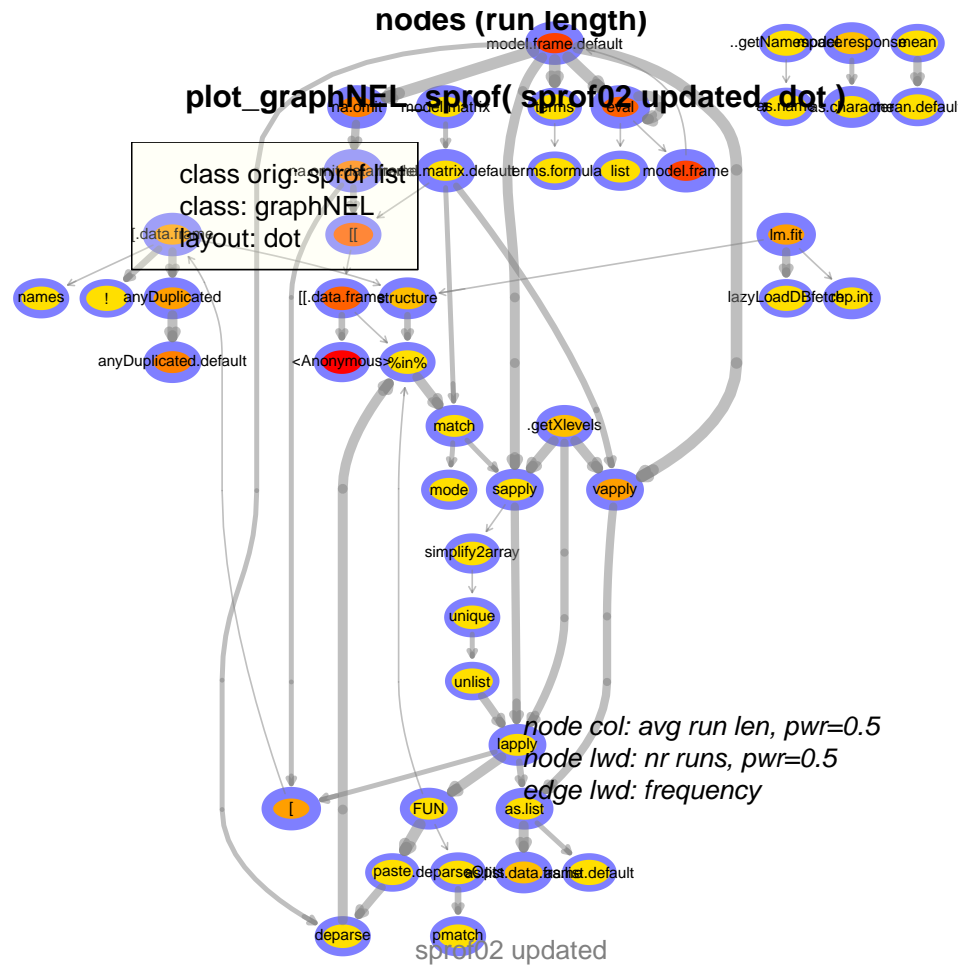
FIGURE 19. nodes (run length)

Input

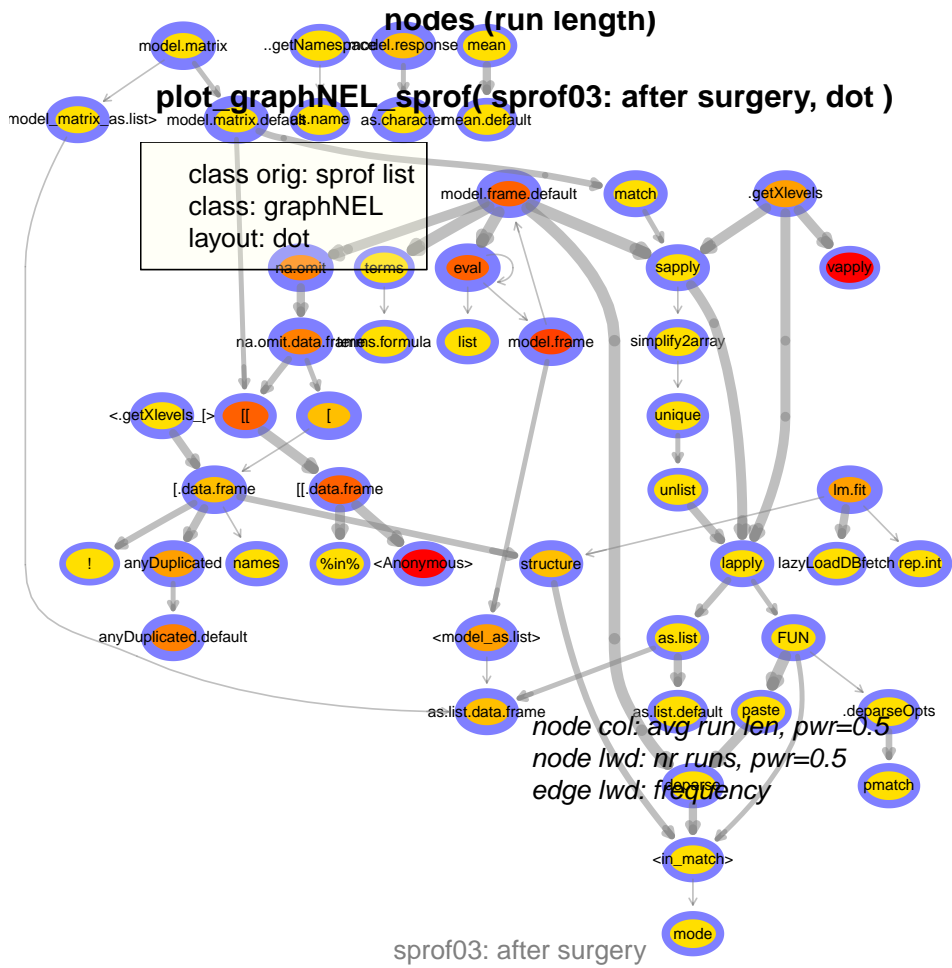
```
plot_graphNEL_sprof(sprof01, nodeattrs="runs")
```



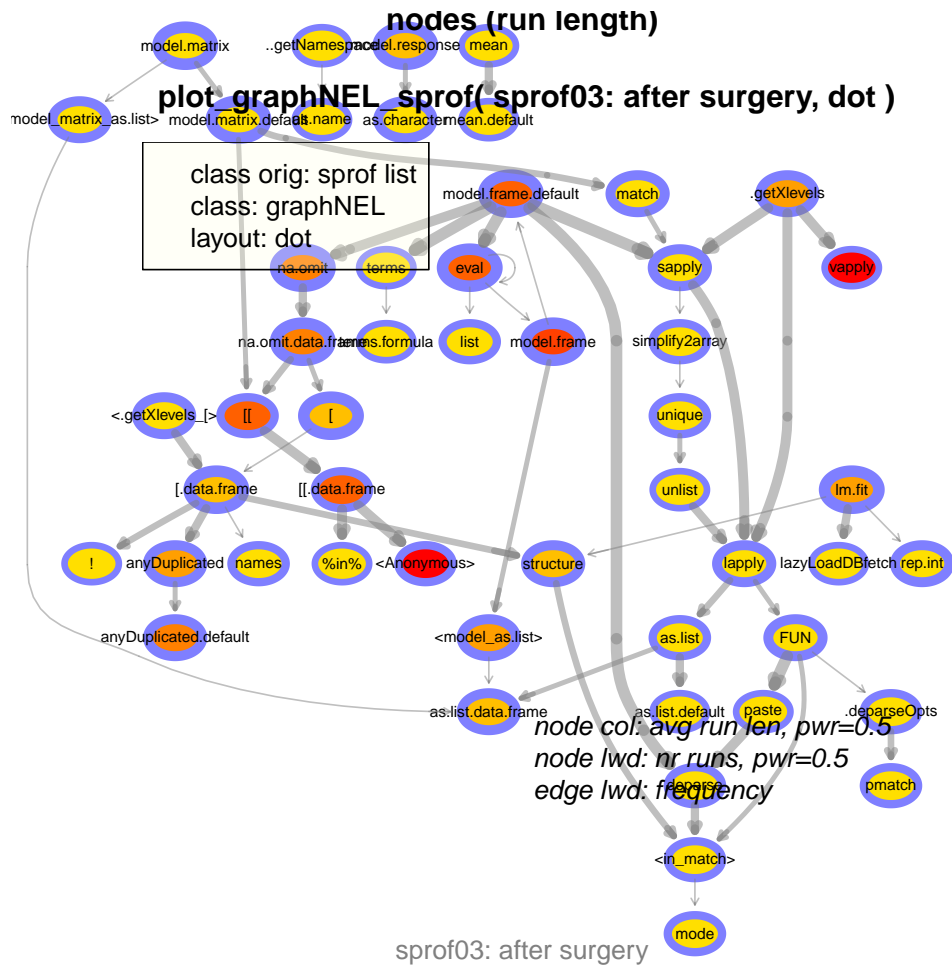
Input
`plot_graphNEL_sprof(sprof02, nodeattrs="runs")`



Input



Input



Now it is time for detailed inspection. We just give hints how to start. Our example is very simple. In our example, we just use two functions, `lm` and `summary`. We should not be surprised to find them at prominent position with dominating weight. The `lm` and `summary.lm` are clearly separated and can be inspected separately.

In the `lm` branch, `eval` is prominent. In the overall structure of R, `eval` is the central interpreter. So it is expected to be prominent in all records.

Sorry. This example may be outdated at the time you read it. This is from R 3.0.1. Chances are that it will be outdated soon when the R core people read it. 2013-08-31

`.getXlevels` is surprising. We have a plain vanilla real numbers problem. Why should it occur at all, and why at a prominent place? It is inside `lm`. But since R is open source, we can inspect it. We find `.getXlevels` at exactly one place near the end of the source of `lm`:

```
z$xlevels <- .getXlevels(mt, mf)
```

`z` is returned as result, and the help file tells us:

`xlevels` (only where relevant) a record of the levels of the factors used in fitting.

In a pure regression problem, it is not relevant. But there is no conditional code at this place. `xlevels` is always calculated (at the cost of some time, and always returning a named list of length 0). Ok. So we found a point where someone could look for an improvement.

We leave additional steps as an exercise. For example: look at the handling of NAs. Why (and where) do we spend time to handle NAs in a problem where there are no NAs at all?

4. STANDARD OUTPUT

For a reference, here are the standard functions.

```
sprof <- sprof01
```

4.1. **Print.** We omit the (lengthy) print output here and just give the commands as a reference.

```
print_nodes(sprof)
```

```
print_stacks(sprof)
```

```
print_profiles(sprof)
```

The `print()` method for *sprof* objects concatenates these three functions.

4.2. **Summary.**

```
summary_nodes(sprof)
```

```
summary_stacks(sprof)
```

```
summary_profiles(sprof)
```

The `summary()` method for *sprof* objects concatenates these three functions.

ToDo: Clarify: "print prints its argument and returns it invisibly (via `invisible(x)`)."

Return the argument, or some print representation?

ToDo: is there a `print=FALSE` variant to postpone printing to e.g. `xtable`?

4.3. Plot. Examples are given in the reference manual for *sprof*.

plot_nodes(sprof) Input

plot_stacks(sprof) Input

plot_profiles(sprof) Input

The `plot.sprof()` method for *sprof* objects concatenates these three functions, see fig. 20. Using the plot functions above allows better control and will be preferred. `shownodes()` may be a sufficient summary, see fig. 10 on page 29.

#12 16
oldpar <- par(mfrow=c(3,4))
plot.sprof(sprof03)
par(oldpar)

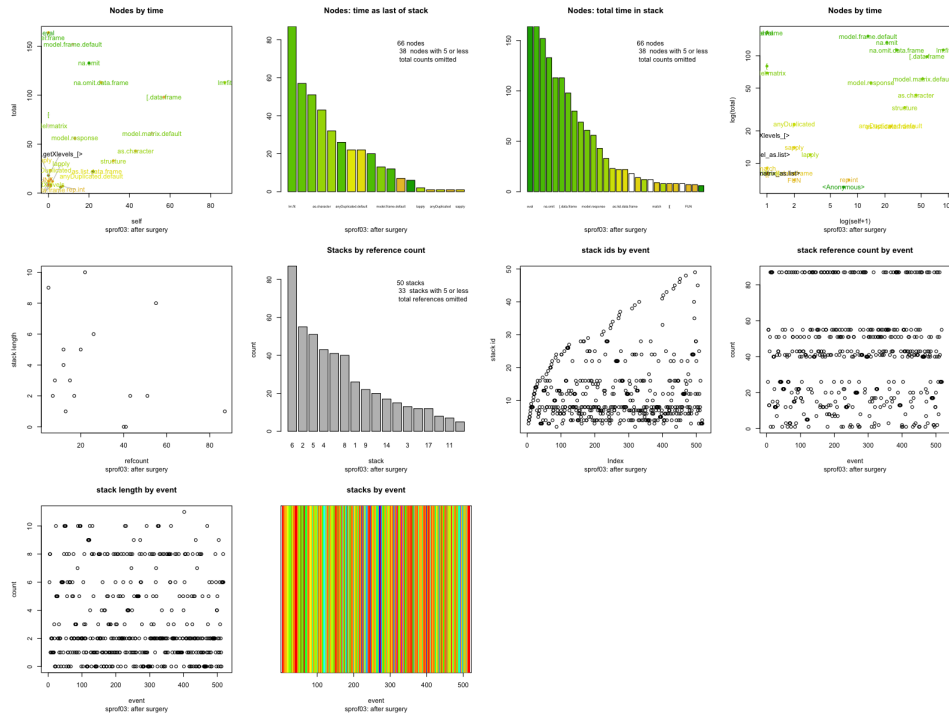


FIGURE 20. `plot.sprof(sprof03)`

5. MORE GRAPHS

Note: This section is collecting experiments with various graph packages. So far, none of the experiments looks too promising. This section is only of interest for you, if you have a preference for some graphic package and want to look up a wrapper here. On the other side: contributions and suggestions are welcome.

Graph layout is a theme of its own. Proposals are readily available, as are their implementation. For some of them, there are R interfaces or re- implementations in R. Their usefulness in our context has to be explored, and the answers will vary with personal preferences.

For some graph layout packages we illustrate an interface here and show a sample result. We use the original profile data here. This is a nasty graph with some R stack peculiarities.

5.1. Example: regression. In this section, we use the recent version of our example, *sprof03* for demonstration. You can re-run it, using your *sprof* data by modifying this instruction by replacing *sprof03* in the next statement with your profile information.

```
sprof <- sprof03
```

To interface *sprof* to a graph handling package, *adjacency()* can extract the adjacency matrix from the profile.

There are various packages for finding a graph layout, and the choice is open to your preferences. The R packages for most of these are just wrapper

```
sprofadj <- adjacency(sprof)
```

This is a format any graph package can handle (maybe). To be on the save side, we provide an (extended) edge list. The added component *lwd* is a proposal for the line width in the graph rendering.

```
sprofedgel <- edgedf(sprofadj)
sprofedgel$lwd <- rkindex(sprofedgel$count,
  maxindex=12,
  ties.method="min")
```

ToDo: by graph package: preferred input format?

ToDo: use attributes. Edge width should be easy.

ToDo: include information from stack connectivity.

ToDo: add usage of *sprofedgel*

This package is now available from Bioconductor only, see <http://www.bioconductor.org/packages/release/bioc/html/graph.html>.

Input

5.1.2. *igraph* package. Attributes for *igraph* are documented in `help(igraph.plotting)`.

Layouts available for *igraph*:

```
layout.auto(graph, dim=2, ...)
layout.random(graph, params, dim=2)
layout.circle(graph, params)
layout.sphere(graph, params)
layout.fruchterman.reingold(graph, ..., dim=2, params)
layout.kamada.kawai(graph, ..., dim=2, params)
layout.spring(graph, ..., params)
layout.reingold.tilford(graph, ..., params)
layout.fruchterman.reingold.grid(graph, ..., params)
layout.lgl(graph, ..., params)
layout.graphopt(graph, ..., params=list())
layout.svd(graph, d=shortest.paths(graph), ...)
layout.norm(layout, xmin = NULL, xmax = NULL, ymin = NULL, ymax = NULL,
zmin = NULL, zmax = NULL) The output of igraph gives problems when rendered
with Adobe Acrobat on OS X Maverick. These examples have been moved to the
demo section. You can call them using
    demo(topic=<Item>, package = "sprof").
```

```
spdemo <- demo("sprof")
prxt(spdemo$results[, c("Item", "Title")])
```

	Item	Title
1	igraphFrRein	igraph Fruchtermann-Reingold layout
2	igraphFrReingrid	igraph Fruchtermann-Reingold grid layout
3	igraphKK	igraph kamada kawai layout
4	igraphKK2	igraph kamada kawai layout alternate
5	igraphauto	igraph auto layout
6	igraphcircle	igraph circle layout
7	igraphgraphopt	igraph graphopt layout
8	igraphlgl	igraph lgl layout
9	igraphrandom	igraph random layout
10	igraphsphere	igraph sphere layout
11	igraphspring	igraph spring layout
12	igraphsvd	igraph svd layout

ToDo: propagate

5.1.3. *network* package.

```

library(network)
as_network_sprof <- function(sprof) {
  sprofadj <- adjacency(sprof)
  nwsprof <- as.network(sprofadj)
  network.vertex.names(nwsprof) <-
    rownames(sprofadj) # not honoured by plot
  return(nwsprof)
}

```

```

plot_network_sprof <- function( nwsprof,
  mode = "fruchtermanreingold",
  main=NULL, label=NULL, sub=NULL,...)
{
  classnwsprof <- class(nwsprof)
  xid <- deparse(substitute(nwsprof))
  xsubid <- NULL

  if (inherits(nwsprof, "sprof")) {
    xsubid <- nwsprof$info$id
    nwsprof <- as_network_sprof(nwsprof)
  }

  if (!is.null(label)) {
    warning("plot_network_sprof: explicit label supplied, but will use vertex names")
    if (!identical(label,network.vertex.names(nwsprof))) {
      print(strx(label))
      print(strx(network.vertex.names(nwsprof)))}
  }

  if (!is.null(sub)) sub <- as.character(sub)
  main <-
    paste0(main,
      "\n plot_network_sprof( ", xid ,", ", mode, " )\n",
      xsubid)

  plot( nwsprof,
    label = network.vertex.names(nwsprof),
    main= main,
    mode = mode,
    edge.len=2,
    edge.col="#80808080",
    sub= sub,
    cex.main=1.5,...)

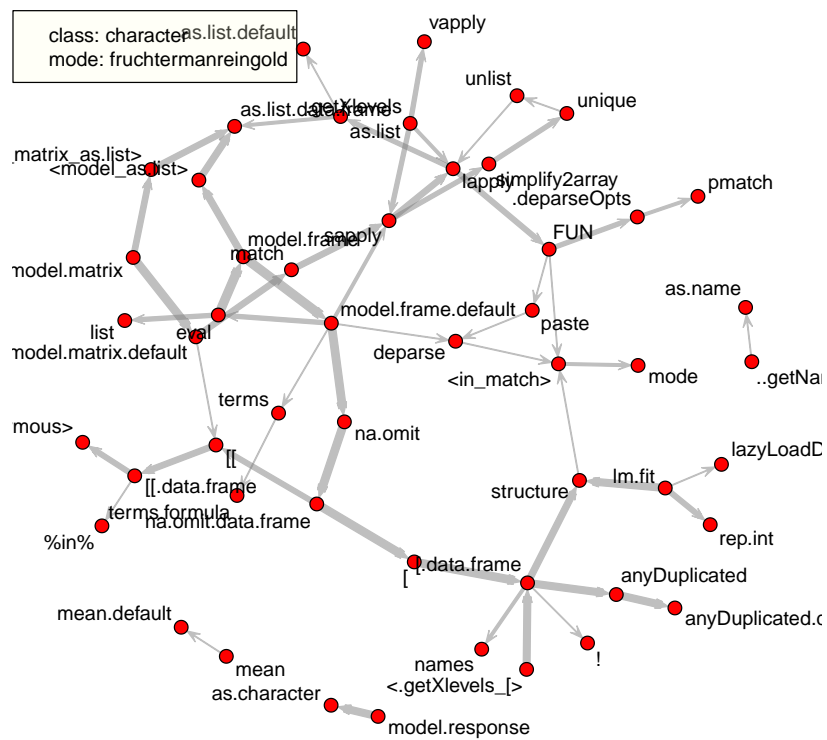
  legend("topleft",
    legend=c(
      paste0("class: ",paste(class(classnwsprof), collapse=" ")),
      paste0("mode: ",mode)),
    bg="#FFFFE040",

```

```
    seg.len=0  
  )  
}
```

```
#8 8
edge.lwd<- sprofadj
edge.lwd[edge.lwd>0]<- rkindex(edge.lwd[edge.lwd>0],
    maxindex=12, ties.method="min")
nwsprof03 <- as.network(sprofadj) # names is not imported
plot_network_sprof(nwsprof03, label=rownames(sprofadj),
    edge.lwd=edge.lwd)
```

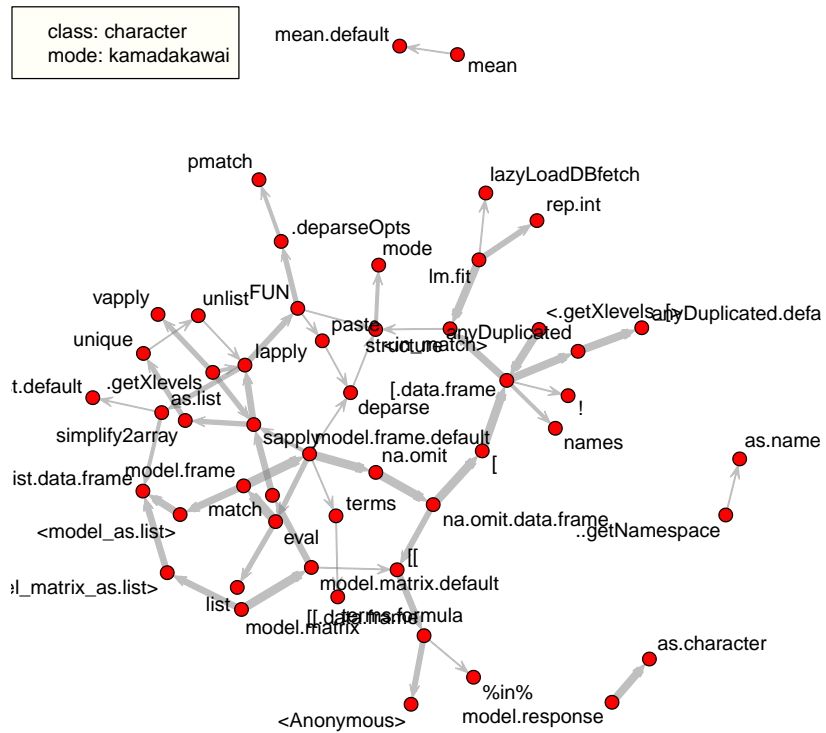
```
plot_network_sprof( nwsprof03, fruchtermanreingold )
```



```
plot_network_sprof(nwsprof03,
  mode="circle",
  label=rownames(sprofadj), edge.lwd=edge.lwd)
```

```
plot_network_sprof(nwsprof03,
  mode="kamadakawai",
  label=rownames(sprofadj), edge.lwd=edge.lwd)
```

```
plot_network_sprof( nwsprof03, kamadakawai )
```



5.1.4. *Rgraphviz package.* Package ‘Rgraphviz’ was removed from the CRAN repository.

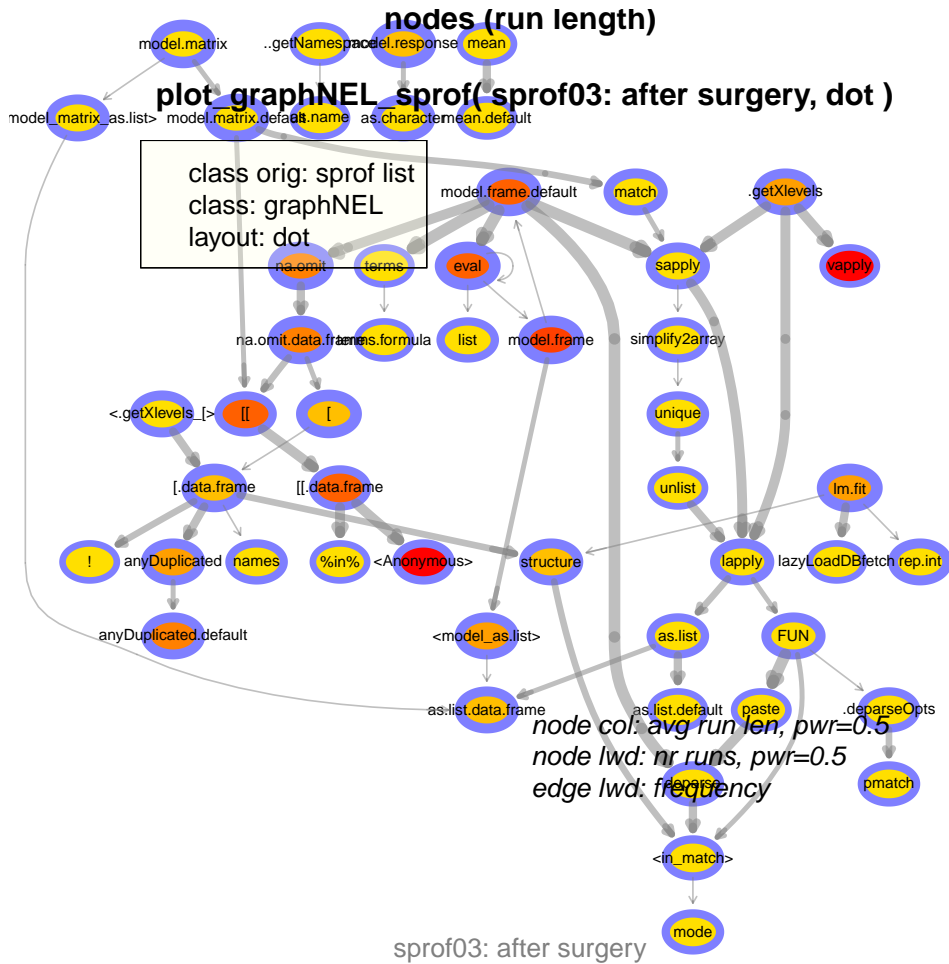
This package is now available from Bioconductor only, see <http://www.bioconductor.org/packages/release/bioc/html/Rgraphviz.html>.

library(Rgraphviz) *Input*

 Input

6 6

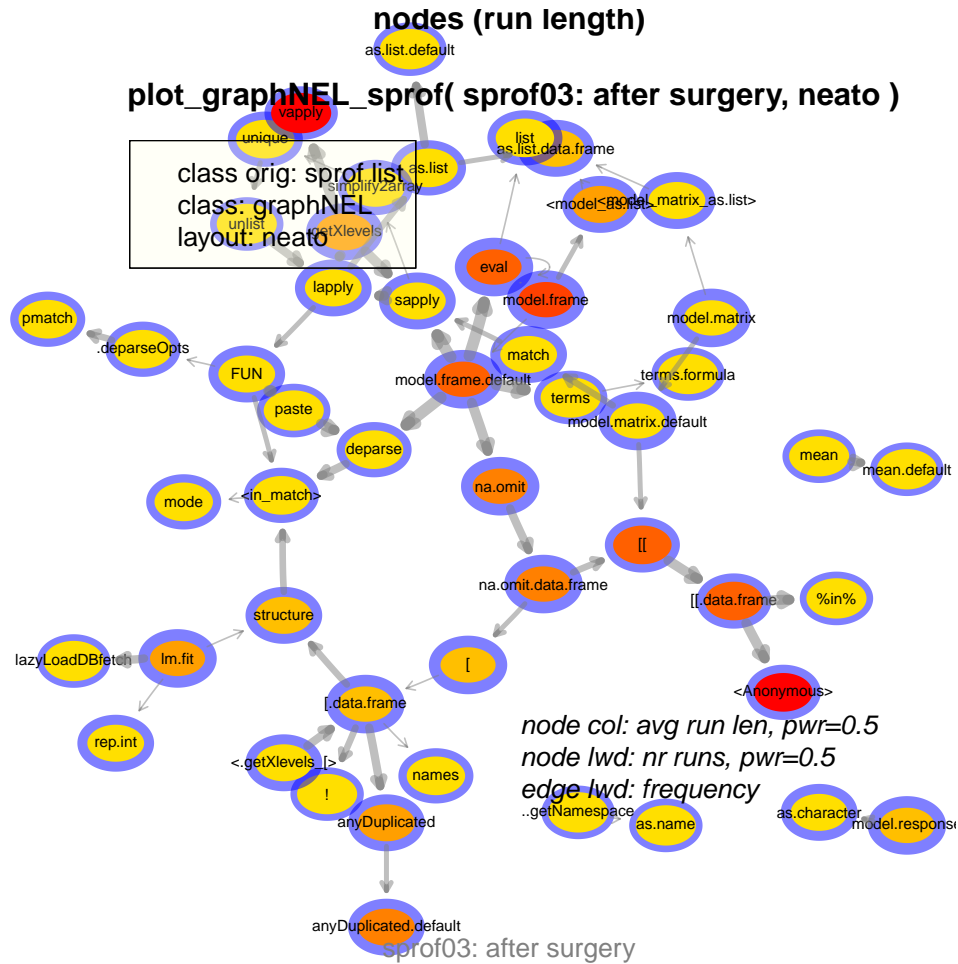
```
plot_graphNEL_sprof(sprof03, layoutType="dot", nodeattrs="runs")
```



Input

#6 6

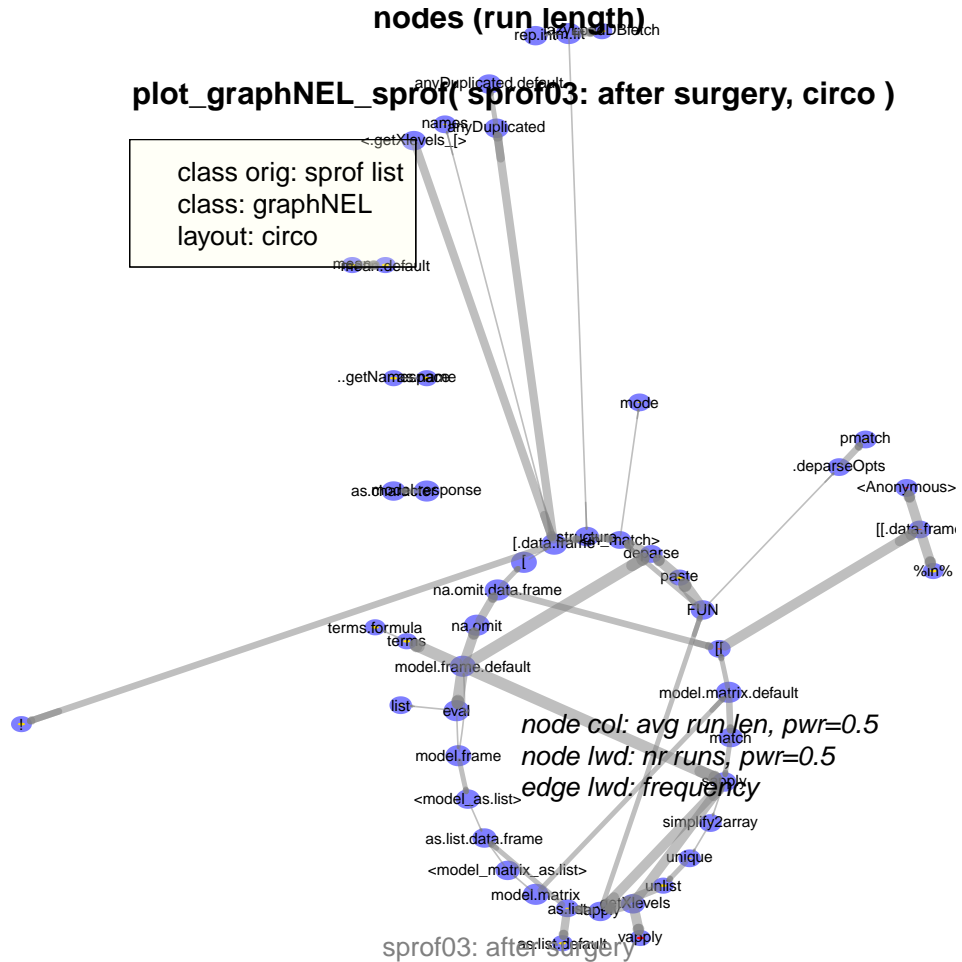
```
plot_graphNEL_sprof(sprof03,
  layoutType="neato", nodeattrs="runs", sub=sprof$info$id)
```



Input

#6 6

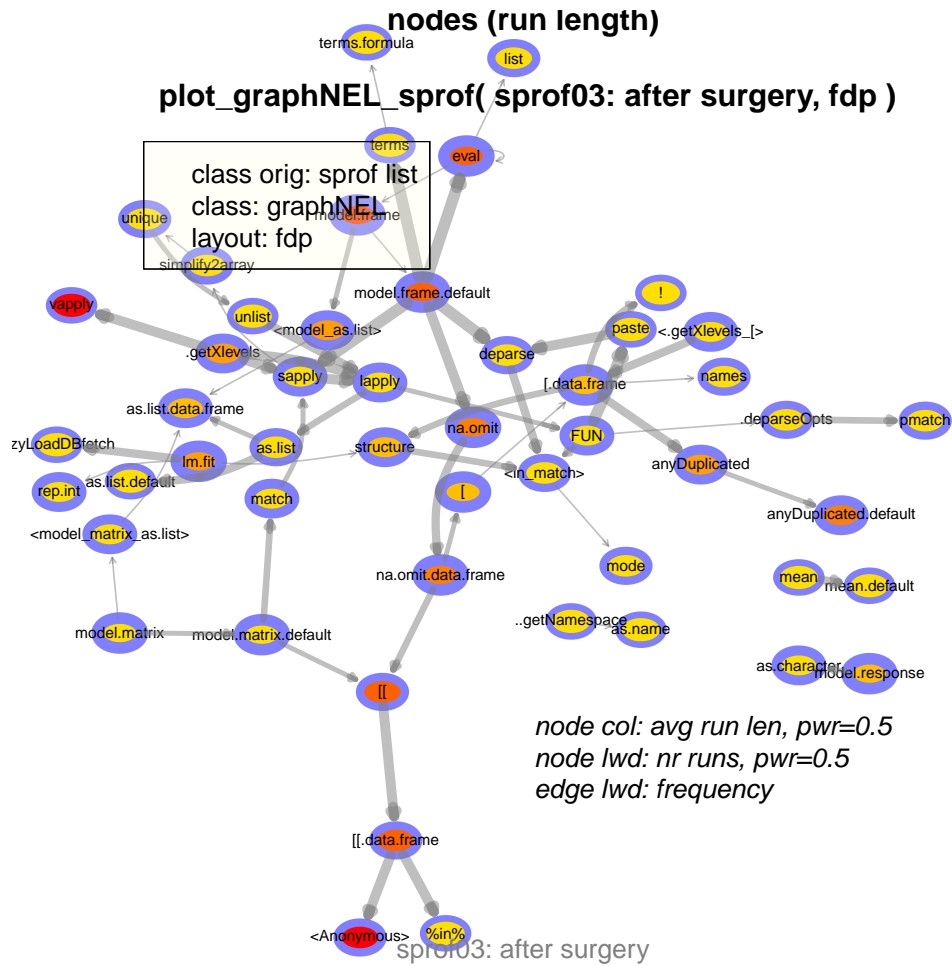
```
plot_graphNEL_sprof(sprof03,
  layoutType="circo", nodeattrs="runs", sub=sprof$info$id)
```



Input

#6 6

```
plot_graphNEL_sprof(sprof03,
  layoutType="fdp", nodeattrs="runs", sub=sprof$info$id)
```



6. TEMPLATE

- Run a profiling routine to profile your functions. You can do it on the fly.
- Read in the profile
- Get a survey
- Trim base level and burn-in/fade-out
- Get a revised survey
- Use a graph display
- Think!

INDEX

Topic **hplot**

plot.sprof, 62
 plot_nodes, 16, 19
 shownodes, 62

Topic **manip**

adjacency, 28
 list.as.matrix, 28
 profiles_matrix, 28
 rrle, 41
 stacks_matrix, 28
 stackstoadj, 28

Topic **misc**

apply, 7
 print, 61
 rkindex, 18, 19
 sampleRprof, 13
 summary, 61
 summaryRprof, 7

Topic **util**

adjacency, 63

ToDo

- 0: remove text vdots from string/name columns. Note: this may be a factor. Use empty string., 5
- 1: Can we calibrate times to CPU rate? Introduce cpu clock cycle as a time base, 8
- 1: Defaults by class, 24
- 1: add as additional information to classical basics, 21
- 1: add class by keyword, 22
- 1: add general discussion of colours in *sprof*, 19
- 1: apply colour to selection?, 19
- 1: check/fix colour for wordcloud, 24
- 1: classes need separate colour palette, distinct from package or keyword., 24
- 1: colour by class – redo. Bundle colour index with colour?, 22
- 1: fix colour. Select colours >after< sorting. Explicit override., 9
- 1: improve barplot.s. Allow vars from matrix or data frame, keep names. Use horizontal names for horizontal layout., 9
- 1: improve colour: support colour in a structure, 20
- 1: make call by reference function, 18
- 1: propagate nodeattrs, 19
- 1: rearrange stacks? detect order?, 7
- 1: remove text vdots from string/name columns, 16
- 1: variable sampling intervals will not be supported in future versions, 13

- 2: stackssrc, collstacksdictrv etc. now out of date, 40
- 2: Implement. Currently best handled on source=text level, 39
- 2: This section needs to be reworked, 38
- 2: Warning: data structure still under discussion, 43
- 2: add a purge function, 34
- 2: add attributes to stacks, and discuss scope, 26
- 2: add marginals and conditionals. Provide function node.summary., 44
- 2: add smart surgery with memory for attributing resources., 41
- 2: add smoothing, remove orphans, 41
- 2: add summary for NA, 44
- 2: check and stabilise colour linking, 28
- 2: check and synchronise, 30
- 2: clean up factor handling, 39
- 2: colours. recolour. Propagate colour to graph., 32
- 2: complete matrix conversion, 28
- 2: cut next level, 39
- 2: data frame conversion?, 28
- 2: fix null name, 41
- 2: function addnode using “call by reference” to be added, 39
- 2: hack. keep length in nodesrunlength, 45
- 2: hack. replace by decent vector/array based implementation, 43
- 2: handle empty stacks and zero counts gracefully, 34
- 2: keep as factor. This is a sparse cube with margins node, stack level, run length. Nodes are mostly concentrated on few levels., 43
- 2: log scale?, 33
- 2: move to other section, 38
- 2: re-think: sort stacks, 26
- 2: replace by sprof03, 41
- 2: should this be NA or NULL?, 34
- 2: sorting/arranging stacks, 26
- 2: table: node #runs min median run length max, 46
- 2: updateRprof needs careful checking. For now, we are including long listings here to provide the necessary information, 33
- 2: warn about undefined vars, e.g. rle, class, ..., 41
- 2: we could do smart smoothing of the stacks here, 31
- 2: xreplacenodes: improve implement, 39
- 3: cut top levels, 48
- 3: merge with as_graphNEL_sprof, 52

- 3: remove global colour; implement local colour, 52
- 4: Clarify: print prints its argument and returns it invisibly (via invisible(x)). Return the argument, or some print representation?, 61
- 4: is there a print=FALSE variant to postpone printing to e.g. xtable?, 61
- 5: add usage of sprofedgel, 63
- 5: by graph package: preferred input format?, 63
- 5: include information from stack connectivity., 63
- 5: propagate, 66
- 5: use attributes. Edge width should be easy., 63

adjacency, 28, 63

apply, 7

Index01, 77

list.as.matrix, 28

node

- attributes, 18

plot.sprof, 62

plot_nodes, 16, 19

print, 61

profiles_matrix, 28

rkindex, 18, 19

rrle, 41

sampleRprof, 13

shownodes, 62

stacks_matrix, 28

stackstoadj, 28

summary, 61

summaryRprof, 7

R session info:

- R version 3.0.2 (2013-09-25), x86_64-apple-darwin10.8.0
- Locale: en_GB.UTF-8/en_GB.UTF-8/en_GB.UTF-8/C/en_GB.UTF-8/en_GB.UTF-8
- Base packages: base, datasets, graphics, grDevices, grid, methods, stats, utils
- Other packages: graph 1.38.3, RColorBrewer 1.0-5, Rcpp 0.10.5, Rgraphviz 2.4.0, sprof 0.1-0, wordcloud 2.4, xtable 1.7-1
- Loaded via a namespace (and not attached): BiocGenerics 0.6.0, network 1.7.2, parallel 3.0.2, slam 0.1-28, sna 2.3-1, stats4 3.0.2, tools 3.0.2

L^AT_EX information:

textwidth: 4.9823in linewidth:4.9823in
 textheight: 8.0824in

Svn repository information:

\$HeadURL: svnssh://gsawitzki@svn.r-forge.r-project.org/svnroot/sintro/pkg/sprof/vignettes/sprofiling.Rnw +
 \$Source: /u/math/j40/cvsroot/lectures/src/insider/profile/Rnw/profile.Rnw,v \$
 \$Id: sprofiling.Rnw 245 2013-10-14 20:15:24Z gsawitzki \$
 \$Revision: 245 \$
 \$Date: 2013-10-14 22:15:24 0200(Mon, 14Oct2013)+
 \$Name: \$
 \$Author: gsawitzki \$

7. XXX – LOST & FOUND

On the profiles level, we know the sample interval length, and the id of the stack recorded. On the stack level, for each stack we have a reference count, with the sample interval lengths used as weights. This reference count is added up for each node in the stack to give the node timings.

Cheap statistics are collected as they come by. For example, from the stacks table it is cheap to identify root and leaf nodes, and this mark is propagated to the nodes table. These are some attempts to recover the factor structures.

GÜNTHER SAWITZKI
STATLAB HEIDELBERG
IM NEUENHEIMER FELD 294
D 69120 HEIDELBERG

E-mail address: `gs@statlab.uni-heidelberg.de`

URL: `http://sintro.r-forge.r-project.org/`