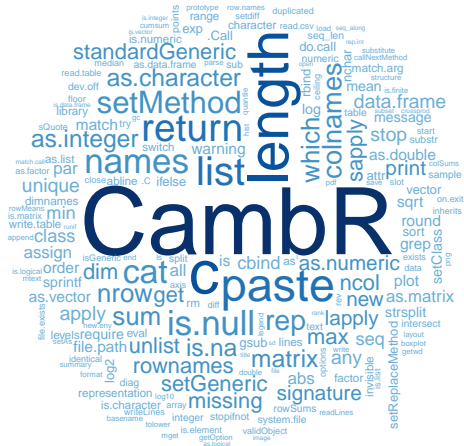


The CambR logo

Laurent Gatto and Robert Stojnic

October 29, 2012

The logo



Motivation

Material for our *Advances R programming course*.

A colourful slide that says

This is what you are expected to know for this course.

Material

- ▶ Get all the code from the **Bioconductor** project

```
svn co https://hedgehog.fhcrc.org/bioconductor/trunk/madman/Rp
```

- ▶ Extract only the .R and .r files

```
find -name "*.rR" | xargs cat > allR.R
```

```
$ ls -sh allR.R
```

```
36M allR.R
```

```
$ wc allR.R
```

```
1008315  3351122 37166763 allR.R
```

Methods - extracting relevant *words*

```
regexp <- "[a-zA-Z.][a-zA-Z0-9._]* *\\(\"
gregexpr(regexp, c("foo", "c (i,j,k)",
                    "setMethod()", "## comment"))
gregexpr(regexp, "foo = c (i,j,k); bar = c(l, m)")

t <- readLines("allR.R")
matches <- gregexpr(regexp, t)
length(matches) ## 1008501
k <- which(sapply(matches, function(x) x[[1]] != -1))
length(k)      ## 502941
```

Methods - counting words

Trim each word by remove leading/ending \t, \n, \f, \r, \s

```
sub("^[\t\n\f\r ]*", "", word)
sub("^[\t\n\f\r ]*$", "", word)
```

Count/increment the word count if is.function(word)

```
if (is.function(word)) {
  if (!(word %in% names(words))) {
    words[[word]] <- 1
  } else {
    words[[word]] <- words[[word]] + 1
  }
}
```

(with some error catching not shown here.)

The output

is a fun/freq data.frame

```
> head(out)
      fun  freq
2       c 38336
11  length 33491
100  paste 22251
17    list 15721
25  return 15236
26    stop 15041
...
```

that needs a bit of post-processing...

Post-processing

- ▶ Take the `sqrt(freq)`
- ▶ Get rid the embarrassing high-freq function stop.
- ▶ Add `CambR` with a `sqrt(freq)` of 300.

Plotting

```
library(wordcloud)
library(RColorBrewer)
pal <- brewer.pal(9,"Blues")[5:9]
wordcloud(out$fun, out$freq, c(6,.1), max.words = 200,
          random.order = FALSE, colors = pal)
```