

Performance

Eric Archer

2/12/2018

Resources

“Optimizing code”: [<http://adv-r.had.co.nz/Profiling.html>]

“A Guide to Speeding Up R Code for Busy People”: [<http://www.noamross.net/blog/2013/4/25/faster-talk.html>]

“Vectorization in R: Why?”: [<http://www.noamross.net/blog/2014/4/16/vectorization-in-r--why.html>]

Performance

In writing code, you are constantly balancing several things. The concept of “performance” could be related to accuracy, speed, stability, readability, or extensibility. In truth, all of these are important, but I would argue that accuracy (getting the right result) comes first. If the code does not deliver what is expected, then it is of no use regardless of how fast it executes or how easy it is to read.

Accuracy is often a direct result of readability. If you write code that is easy to follow and read, the flow of your logic will be obvious and tend to be clear - most importantly, to you!

Even when you’ve written some code that is performing as expected, you may want to spend some time getting it to run faster. The very first thing you need to do is understand where the bottlenecks are. Depending on how you’ve structured your code, there are multiple things that could be slower than expected and can be sped up with some restructuring.

Note that everything takes some time, so nothing is truly instantaneous. Recognizing that means that you will only be able to optimize but so much. At some point you will be spending more time re-programming trying to eke out a handful of milliseconds than it would take for all of the runs you could possibly imagine. Stop. You’ve done enough.

Profiling

But, if you’re just getting started in this process, the first step you should do is profile your code. The core R functions for this are `Rprof` and `summaryRprof`. The former is used to start and stop profiling of code that has executed, and the latter summarizes the results of that profiling. In the example below, we’ll profile some code that does some permutation of CTD data to produce a bootstrap mean of the temperature at each station.

```
# we're doing 10 replicates and want to return the result in an array
boot.mean <- sapply(1:10, function(i) {
  # read the data in
  df <- read.csv("ctd.csv", stringsAsFactors = FALSE)
  # loop over each value of station
  sapply(unique(df$station), function(st) {
    # identify the rows for this station and the first depth level
    i <- which(df$station == st & df$depth == 1)
    # take a random sample of these rows (with replacement)
    i <- sample(i, length(i), replace = TRUE)
```

```

    # return the mean of temperature for these rows
    mean(df$temp[i])
  })
})
ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))
ci

```

	Station.1	Station.10	Station.11	Station.12	Station.13	Station.14
2.5%	16.65381	16.19743	15.99068	16.43686	16.60606	16.66461
97.5%	17.79979	16.93009	16.67294	17.17486	17.70802	17.42815
	Station.15	Station.16	Station.17	Station.18	Station.19	Station.2
2.5%	16.20301	16.17403	16.18031	15.86168	16.92970	16.44211
97.5%	17.19676	17.31640	16.92154	16.60077	17.41356	17.30509
	Station.20	Station.21	Station.22	Station.23	Station.24	Station.25
2.5%	16.53933	16.66187	16.37522	16.40888	17.26842	16.69136
97.5%	17.28435	17.56595	17.63911	16.99161	17.73762	16.92737
	Station.26	Station.27	Station.28	Station.29	Station.3	Station.30
2.5%	16.74014	15.83145	16.54754	16.0018	16.30871	16.04785
97.5%	16.98671	16.65807	17.71411	17.1856	17.00311	16.92396
	Station.31	Station.32	Station.33	Station.34	Station.35	Station.36
2.5%	16.23272	17.12982	15.88091	16.13300	16.29069	16.11228
97.5%	16.79729	17.80915	16.66506	16.52338	16.79769	17.21335
	Station.37	Station.38	Station.39	Station.4	Station.40	Station.5
2.5%	15.70534	16.65742	16.84594	15.94330	16.98624	15.44500
97.5%	16.42465	17.15512	17.03477	16.67976	17.84088	16.29437
	Station.6	Station.7	Station.8	Station.9		
2.5%	16.07635	16.61468	16.46316	16.54747		
97.5%	17.05911	17.55426	17.56230	17.27378		

Let's first see what parts of this are taking the most time and how much:

```

# Start profiling
Rprof()

# Run code
boot.mean <- sapply(1:10, function(i) {
  df <- read.csv("ctd.csv", stringsAsFactors = FALSE)
  sapply(unique(df$station), function(st) {
    i <- which(df$station == st & df$depth == 1)
    i <- sample(i, length(i), replace = TRUE)
    mean(df$temp[i])
  })
})

ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))

# Stop profiling
Rprof(NULL)

# Examine profile summary
summaryRprof()

```

\$by.self	self.time	self.pct	total.time	total.pct
"scan"	5.70	86.36	5.70	86.36
"which"	0.50	7.58	0.52	7.88

```

".External2"          0.34    5.15    0.34    5.15
"[.data.frame"        0.02    0.30    0.02    0.30
"close.connection"    0.02    0.30    0.02    0.30
"unique.default"      0.02    0.30    0.02    0.30

$by.total
      total.time total.pct self.time self.pct
"block_exec"      6.60   100.00    0.00    0.00
"call_block"      6.60   100.00    0.00    0.00
"doTryCatch"      6.60   100.00    0.00    0.00
"eval"            6.60   100.00    0.00    0.00
"evaluate_call"   6.60   100.00    0.00    0.00
"evaluate::evaluate" 6.60   100.00    0.00    0.00
"evaluate"        6.60   100.00    0.00    0.00
"FUN"             6.60   100.00    0.00    0.00
"handle"          6.60   100.00    0.00    0.00
"in_dir"          6.60   100.00    0.00    0.00
"knitr::knit"     6.60   100.00    0.00    0.00
"lapply"          6.60   100.00    0.00    0.00
"process_file"    6.60   100.00    0.00    0.00
"process_group.block" 6.60   100.00    0.00    0.00
"process_group"   6.60   100.00    0.00    0.00
"rmarkdown::render" 6.60   100.00    0.00    0.00
"sapply"          6.60   100.00    0.00    0.00
"timing_fn"       6.60   100.00    0.00    0.00
"try"            6.60   100.00    0.00    0.00
"tryCatch"        6.60   100.00    0.00    0.00
"tryCatchList"    6.60   100.00    0.00    0.00
"tryCatchOne"     6.60   100.00    0.00    0.00
"withCallingHandlers" 6.60   100.00    0.00    0.00
"withVisible"     6.60   100.00    0.00    0.00
"read.csv"        6.06    91.82    0.00    0.00
"read.table"      6.06    91.82    0.00    0.00
"scan"           5.70    86.36    5.70    86.36
"which"           0.52     7.88    0.50     7.58
".External2"      0.34     5.15    0.34     5.15
"type.convert"    0.34     5.15    0.00     0.00
"[.data.frame"    0.02     0.30    0.02     0.30
"close.connection" 0.02     0.30    0.02     0.30
"unique.default"  0.02     0.30    0.02     0.30
"["              0.02     0.30    0.00     0.00
"$data.frame"     0.02     0.30    0.00     0.00
"$"              0.02     0.30    0.00     0.00
"close"           0.02     0.30    0.00     0.00
"unique"          0.02     0.30    0.00     0.00

$sample.interval
[1] 0.02

$sampling.time
[1] 6.6

```

We can see that a majority of the time was taken by reading the file (“scan”, “read.table”, “read.csv”). This makes sense because it is doing that every iteration of the loop in the `sapply` function. Let’s be smarter and

change the code so that file reading is happening only once:

```
# Start profiling
Rprof()

# Run code
df <- read.csv("ctd.csv", stringsAsFactors = FALSE)
boot.mean <- sapply(1:10, function(i) {
  sapply(unique(df$station), function(st) {
    i <- which(df$station == st & df$depth == 1)
    i <- sample(i, length(i), replace = TRUE)
    mean(df$temp[i])
  })
})
ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))

# Stop profiling
Rprof(NULL)

# Examine profile summary
summaryRprof()
```

```
$by.self
      self.time self.pct total.time total.pct
"scan"      0.60   53.57      0.60   53.57
"which"      0.44   39.29      0.48   42.86
".External2" 0.04    3.57      0.04    3.57
"&"          0.02    1.79      0.02    1.79
"=="         0.02    1.79      0.02    1.79

$by.total
      total.time total.pct self.time self.pct
"block_exec"      1.12  100.00      0.00    0.00
"call_block"      1.12  100.00      0.00    0.00
"doTryCatch"      1.12  100.00      0.00    0.00
"eval"            1.12  100.00      0.00    0.00
"evaluate_call"   1.12  100.00      0.00    0.00
"evaluate::evaluate" 1.12  100.00      0.00    0.00
"evaluate"        1.12  100.00      0.00    0.00
"handle"          1.12  100.00      0.00    0.00
"in_dir"          1.12  100.00      0.00    0.00
"knitr::knit"     1.12  100.00      0.00    0.00
"process_file"    1.12  100.00      0.00    0.00
"process_group.block" 1.12  100.00      0.00    0.00
"process_group"   1.12  100.00      0.00    0.00
"rmarkdown::render" 1.12  100.00      0.00    0.00
"timing_fn"        1.12  100.00      0.00    0.00
"try"             1.12  100.00      0.00    0.00
"tryCatch"        1.12  100.00      0.00    0.00
"tryCatchList"    1.12  100.00      0.00    0.00
"tryCatchOne"     1.12  100.00      0.00    0.00
"withCallingHandlers" 1.12  100.00      0.00    0.00
"withVisible"     1.12  100.00      0.00    0.00
"read.csv"        0.64   57.14      0.00    0.00
"read.table"      0.64   57.14      0.00    0.00
```

"scan"	0.60	53.57	0.60	53.57
"which"	0.48	42.86	0.44	39.29
"FUN"	0.48	42.86	0.00	0.00
"lapply"	0.48	42.86	0.00	0.00
"sapply"	0.48	42.86	0.00	0.00
".External2"	0.04	3.57	0.04	3.57
"type.convert"	0.04	3.57	0.00	0.00
"&"	0.02	1.79	0.02	1.79
"=="	0.02	1.79	0.02	1.79

```
$sample.interval
[1] 0.02
```

```
$sampling.time
[1] 1.12
```

Notice that our total time decreased, but file reading is still taking the most time. Since this is on its own line, we can't really make this any faster, but let's profile just the nested `sapply` lines:

```
df <- read.csv("ctd.csv", stringsAsFactors = FALSE)

# Start profiling
Rprof()

# Run code
boot.mean <- sapply(1:10, function(i) {
  sapply(unique(df$station), function(st) {
    i <- which(df$station == st & df$depth == 1)
    i <- sample(i, length(i), replace = TRUE)
    mean(df$temp[i])
  })
})
ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))

# Stop profiling
Rprof(NULL)

# Examine profile summary
summaryRprof()
```

```
$by.self
      self.time self.pct total.time total.pct
"which"      0.38   79.17      0.46   95.83
"&"          0.04    8.33      0.04    8.33
"=="         0.04    8.33      0.04    8.33
"is.factor"  0.02    4.17      0.02    4.17
```

```
$by.total
      total.time total.pct self.time self.pct
"block_exec"      0.48  100.00      0.00    0.00
"call_block"      0.48  100.00      0.00    0.00
"doTryCatch"      0.48  100.00      0.00    0.00
"eval"            0.48  100.00      0.00    0.00
"evaluate_call"   0.48  100.00      0.00    0.00
"evaluate::evaluate" 0.48  100.00      0.00    0.00
```

"evaluate"	0.48	100.00	0.00	0.00
"FUN"	0.48	100.00	0.00	0.00
"handle"	0.48	100.00	0.00	0.00
"in_dir"	0.48	100.00	0.00	0.00
"knitr::knit"	0.48	100.00	0.00	0.00
"process_file"	0.48	100.00	0.00	0.00
"process_group.block"	0.48	100.00	0.00	0.00
"process_group"	0.48	100.00	0.00	0.00
"rmarkdown::render"	0.48	100.00	0.00	0.00
"timing_fn"	0.48	100.00	0.00	0.00
"try"	0.48	100.00	0.00	0.00
"tryCatch"	0.48	100.00	0.00	0.00
"tryCatchList"	0.48	100.00	0.00	0.00
"tryCatchOne"	0.48	100.00	0.00	0.00
"withCallingHandlers"	0.48	100.00	0.00	0.00
"withVisible"	0.48	100.00	0.00	0.00
"which"	0.46	95.83	0.38	79.17
"lapply"	0.46	95.83	0.00	0.00
"sapply"	0.46	95.83	0.00	0.00
"&"	0.04	8.33	0.04	8.33
"=="	0.04	8.33	0.04	8.33
"is.factor"	0.02	4.17	0.02	4.17
"apply"	0.02	4.17	0.00	0.00
"quantile.default"	0.02	4.17	0.00	0.00

```
$sample.interval
```

```
[1] 0.02
```

```
$sampling.time
```

```
[1] 0.48
```

Now it seems like `which` is taking most of the time, so let's focus on that. If we remember that we don't have to use `which`. We can just index with the base logical vector. Since we still need to randomly sample temperature with replacement, we'll extract that column and do the sampling on it:

```
df <- read.csv("ctd.csv", stringsAsFactors = FALSE)

# Start profiling
Rprof()

# Run code
boot.mean <- sapply(1:10, function(i) {
  sapply(unique(df$station), function(st) {
    temp <- df$temp[df$station == st & df$depth == 1]
    temp <- sample(temp, length(temp), replace = TRUE)
    mean(temp)
  })
})

ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))

# Stop profiling
Rprof(NULL)

# Examine profile summary
summaryRprof()
```

```

$by.self
      self.time self.pct total.time total.pct
"FUN"      0.46      92      0.50      100
"&"        0.02       4      0.02       4
"unique.default" 0.02       4      0.02       4

$by.total
      total.time total.pct self.time self.pct
"FUN"      0.50      100      0.46      92
"block_exec" 0.50      100      0.00       0
"call_block" 0.50      100      0.00       0
"doTryCatch" 0.50      100      0.00       0
"eval"       0.50      100      0.00       0
"evaluate_call" 0.50      100      0.00       0
"evaluate::evaluate" 0.50      100      0.00       0
"evaluate"   0.50      100      0.00       0
"handle"     0.50      100      0.00       0
"in_dir"     0.50      100      0.00       0
"knitr::knit" 0.50      100      0.00       0
"lapply"     0.50      100      0.00       0
"process_file" 0.50      100      0.00       0
"process_group.block" 0.50      100      0.00       0
"process_group" 0.50      100      0.00       0
"rmarkdown::render" 0.50      100      0.00       0
"sapply"     0.50      100      0.00       0
"timing_fn"   0.50      100      0.00       0
"try"        0.50      100      0.00       0
"tryCatch"   0.50      100      0.00       0
"tryCatchList" 0.50      100      0.00       0
"tryCatchOne" 0.50      100      0.00       0
"withCallingHandlers" 0.50      100      0.00       0
"withVisible" 0.50      100      0.00       0
"&"         0.02       4      0.02       4
"unique.default" 0.02       4      0.02       4
"unique"     0.02       4      0.00       0

```

```

$sample.interval
[1] 0.02

```

```

$sampling.time
[1] 0.5

```

We can now see that we're spending a non-negligible amount of time in the logical operators, == and &. Since we're always interested in the first depth class (depth == 1), let's extract that early on:

```

df <- read.csv("ctd.csv", stringsAsFactors = FALSE)

# Start profiling
Rprof()

# Run code
df.1 <- df[df$depth == 1, ]
boot.mean <- sapply(1:10, function(i) {
  sapply(unique(df.1$station), function(st) {
    temp <- df.1$temp[df.1$station == st]

```

```

    temp <- sample(temp, length(temp), replace = TRUE)
    mean(temp)
  })
})
ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))

# Stop profiling
Rprof(NULL)

# Examine profile summary
summaryRprof()

```

```

$by.self
              self.time self.pct total.time total.pct
"[".data.frame"    0.02    100      0.02      100

$by.total
              total.time total.pct self.time self.pct
"[".data.frame"    0.02    100      0.02      100
"["               0.02    100      0.00        0
"$ .data.frame"    0.02    100      0.00        0
"$               0.02    100      0.00        0
"block_exec"       0.02    100      0.00        0
"call_block"       0.02    100      0.00        0
"doTryCatch"       0.02    100      0.00        0
"eval"             0.02    100      0.00        0
"evaluate_call"    0.02    100      0.00        0
"evaluate::evaluate" 0.02    100      0.00        0
"evaluate"         0.02    100      0.00        0
"FUN"              0.02    100      0.00        0
"handle"           0.02    100      0.00        0
"in_dir"           0.02    100      0.00        0
"knitr::knit"      0.02    100      0.00        0
"lapply"           0.02    100      0.00        0
"process_file"     0.02    100      0.00        0
"process_group.block" 0.02    100      0.00        0
"process_group"    0.02    100      0.00        0
"rmarkdown::render" 0.02    100      0.00        0
"sapply"           0.02    100      0.00        0
"timing_fn"         0.02    100      0.00        0
"try"              0.02    100      0.00        0
"tryCatch"         0.02    100      0.00        0
"tryCatchList"     0.02    100      0.00        0
"tryCatchOne"      0.02    100      0.00        0
"withCallingHandlers" 0.02    100      0.00        0
"withVisible"      0.02    100      0.00        0

$sample.interval
[1] 0.02

$sampling.time
[1] 0.02

```

This is considerably faster than before. Let's see where (if) there are other bottlenecks now by increasing

the number of replicates from 10 to 1000. Also, because the total time is getting smaller, let's decrease the sampling interval to 0.01.

```
df <- read.csv("ctd.csv", stringsAsFactors = FALSE)

# Start profiling
Rprof(interval = 0.01)

# Run code
df.1 <- df[df$depth == 1, ]
boot.mean <- sapply(1:1000, function(i) {
  sapply(unique(df.1$station), function(st) {
    temp <- df.1$temp[df.1$station == st]
    temp <- sample(temp, length(temp), replace = TRUE)
    mean(temp)
  })
})
ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))

# Stop profiling
Rprof(NULL)

# Examine profile summary
summaryRprof()
```

```
$by.self
```

	self.time	self.pct	total.time	total.pct
"FUN"	1.26	67.38	1.86	99.47
"sample"	0.11	5.88	0.23	12.30
"sample.int"	0.11	5.88	0.11	5.88
"["	0.07	3.74	0.16	8.56
"unique.default"	0.07	3.74	0.07	3.74
"mean"	0.05	2.67	0.06	3.21
"\$"	0.03	1.60	0.20	10.70
"[.data.frame"	0.03	1.60	0.09	4.81
"lapply"	0.02	1.07	1.86	99.47
"%in%"	0.02	1.07	0.04	2.14
"length"	0.02	1.07	0.02	1.07
"\$.data.frame"	0.01	0.53	0.17	9.09
"\$<-"	0.01	0.53	0.01	0.53
"all"	0.01	0.53	0.01	0.53
"apply"	0.01	0.53	0.01	0.53
"names"	0.01	0.53	0.01	0.53
"nargs"	0.01	0.53	0.01	0.53
"simplify2array"	0.01	0.53	0.01	0.53
"sys.call"	0.01	0.53	0.01	0.53

```
$by.total
```

	total.time	total.pct	self.time	self.pct
"block_exec"	1.87	100.00	0.00	0.00
"call_block"	1.87	100.00	0.00	0.00
"doTryCatch"	1.87	100.00	0.00	0.00
"eval"	1.87	100.00	0.00	0.00
"evaluate_call"	1.87	100.00	0.00	0.00
"evaluate::evaluate"	1.87	100.00	0.00	0.00

"evaluate"	1.87	100.00	0.00	0.00
"handle"	1.87	100.00	0.00	0.00
"in_dir"	1.87	100.00	0.00	0.00
"knitr::knit"	1.87	100.00	0.00	0.00
"process_file"	1.87	100.00	0.00	0.00
"process_group.block"	1.87	100.00	0.00	0.00
"process_group"	1.87	100.00	0.00	0.00
"rmarkdown::render"	1.87	100.00	0.00	0.00
"timing_fn"	1.87	100.00	0.00	0.00
"try"	1.87	100.00	0.00	0.00
"tryCatch"	1.87	100.00	0.00	0.00
"tryCatchList"	1.87	100.00	0.00	0.00
"tryCatchOne"	1.87	100.00	0.00	0.00
"withCallingHandlers"	1.87	100.00	0.00	0.00
"withVisible"	1.87	100.00	0.00	0.00
"FUN"	1.86	99.47	1.26	67.38
"lapply"	1.86	99.47	0.02	1.07
"sapply"	1.86	99.47	0.00	0.00
"sample"	0.23	12.30	0.11	5.88
"\$"	0.20	10.70	0.03	1.60
"\$.data.frame"	0.17	9.09	0.01	0.53
"[["	0.16	8.56	0.07	3.74
"sample.int"	0.11	5.88	0.11	5.88
"[.data.frame"	0.09	4.81	0.03	1.60
"unique"	0.08	4.28	0.00	0.00
"unique.default"	0.07	3.74	0.07	3.74
"mean"	0.06	3.21	0.05	2.67
"%in%"	0.04	2.14	0.02	1.07
"length"	0.02	1.07	0.02	1.07
"\$<-"	0.01	0.53	0.01	0.53
"all"	0.01	0.53	0.01	0.53
"apply"	0.01	0.53	0.01	0.53
"names"	0.01	0.53	0.01	0.53
"nargs"	0.01	0.53	0.01	0.53
"simplify2array"	0.01	0.53	0.01	0.53
"sys.call"	0.01	0.53	0.01	0.53
"cb\$putconst"	0.01	0.53	0.00	0.00
"cmp"	0.01	0.53	0.00	0.00
"cmpCall"	0.01	0.53	0.00	0.00
"cmpCallArgs"	0.01	0.53	0.00	0.00
"cmpCallSymFun"	0.01	0.53	0.00	0.00
"cmpfun"	0.01	0.53	0.00	0.00
"compiler:::tryCmpfun"	0.01	0.53	0.00	0.00
"genCode"	0.01	0.53	0.00	0.00
"h"	0.01	0.53	0.00	0.00
"make.promiseContext"	0.01	0.53	0.00	0.00
"mean.default"	0.01	0.53	0.00	0.00
"tryInline"	0.01	0.53	0.00	0.00

\$sample.interval

[1] 0.01

\$sampling.time

[1] 1.87

Here we see that we're spending most of our recoverable time in sample. It isn't easy to make that or simplify that step as we have to do this random sampling. The next items down have to do with indexing the data frame (`[.data.frame]`). Lets try to handle that by doing some extraction ahead of time and working with vectors rather than data.frames:

```
df <- read.csv("ctd.csv", stringsAsFactors = FALSE)

# Start profiling
Rprof(interval = 0.01)

# Run code
df.1 <- df[df$depth == 1, ]
temp <- df.1$temp
station <- df.1$station

boot.mean <- sapply(1:1000, function(i) {
  sapply(unique(station), function(st) {
    st.temp <- temp[station == st]
    st.temp <- sample(st.temp, length(st.temp), replace = TRUE)
    mean(st.temp)
  })
})
ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))

# Stop profiling
Rprof(NULL)

# Examine profile summary
summaryRprof()
```

```
$by.self
      self.time self.pct total.time total.pct
"FUN"          1.11   77.62         1.43   100.00
"sample.int"    0.11    7.69          0.11    7.69
"unique.default" 0.06    4.20          0.06    4.20
"mean"          0.05    3.50          0.07    4.90
"sample"        0.04    2.80          0.15   10.49
"lapply"        0.02    1.40          1.43   100.00
"mean.default"  0.02    1.40          0.02    1.40
"any.dots"      0.01    0.70          0.01    0.70
"lengths"       0.01    0.70          0.01    0.70

$by.total
      total.time total.pct self.time self.pct
"FUN"          1.43   100.00         1.11   77.62
"lapply"        1.43   100.00          0.02    1.40
"block_exec"    1.43   100.00          0.00    0.00
"call_block"    1.43   100.00          0.00    0.00
"doTryCatch"    1.43   100.00          0.00    0.00
"eval"          1.43   100.00          0.00    0.00
"evaluate_call"  1.43   100.00          0.00    0.00
"evaluate::evaluate" 1.43   100.00          0.00    0.00
"evaluate"      1.43   100.00          0.00    0.00
"handle"        1.43   100.00          0.00    0.00
"in_dir"        1.43   100.00          0.00    0.00
```

"knitr::knit"	1.43	100.00	0.00	0.00
"process_file"	1.43	100.00	0.00	0.00
"process_group.block"	1.43	100.00	0.00	0.00
"process_group"	1.43	100.00	0.00	0.00
"rmarkdown::render"	1.43	100.00	0.00	0.00
"sapply"	1.43	100.00	0.00	0.00
"timing_fn"	1.43	100.00	0.00	0.00
"try"	1.43	100.00	0.00	0.00
"tryCatch"	1.43	100.00	0.00	0.00
"tryCatchList"	1.43	100.00	0.00	0.00
"tryCatchOne"	1.43	100.00	0.00	0.00
"withCallingHandlers"	1.43	100.00	0.00	0.00
"withVisible"	1.43	100.00	0.00	0.00
"sample"	0.15	10.49	0.04	2.80
"sample.int"	0.11	7.69	0.11	7.69
"mean"	0.07	4.90	0.05	3.50
"unique"	0.07	4.90	0.00	0.00
"unique.default"	0.06	4.20	0.06	4.20
"mean.default"	0.02	1.40	0.02	1.40
"any.dots"	0.01	0.70	0.01	0.70
"lengths"	0.01	0.70	0.01	0.70
"cb\$putconst"	0.01	0.70	0.00	0.00
"checkCall"	0.01	0.70	0.00	0.00
"cmp"	0.01	0.70	0.00	0.00
"cmpCall"	0.01	0.70	0.00	0.00
"cmpCallArgs"	0.01	0.70	0.00	0.00
"cmpCallSymFun"	0.01	0.70	0.00	0.00
"cmpfun"	0.01	0.70	0.00	0.00
"cmpSymbolAssign"	0.01	0.70	0.00	0.00
"compiler::tryCmpfun"	0.01	0.70	0.00	0.00
"genCode"	0.01	0.70	0.00	0.00
"h"	0.01	0.70	0.00	0.00
"simplify2array"	0.01	0.70	0.00	0.00
"tryInline"	0.01	0.70	0.00	0.00

```
$sample.interval
[1] 0.01
```

```
$sampling.time
[1] 1.43
```

That made a noticeable improvement. It also highlights something else we're doing repeatedly - using `unique` to get the unique station names. Lets do that earlier.

```
df <- read.csv("ctd.csv", stringsAsFactors = FALSE)

# Start profiling
Rprof(interval = 0.01)

# Run code
df.1 <- df[df$depth == 1, ]
temp <- df.1$temp
station <- df.1$station
st.names <- unique(station)
```

```

boot.mean <- sapply(1:1000, function(i) {
  sapply(st.names, function(st) {
    st.temp <- temp[station == st]
    st.temp <- sample(st.temp, length(st.temp), replace = TRUE)
    mean(st.temp)
  })
})
ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))

# Stop profiling
Rprof(NULL)

# Examine profile summary
summaryRprof()

```

\$by.self

	self.time	self.pct	total.time	total.pct
"FUN"	1.05	76.64	1.36	99.27
"sample.int"	0.11	8.03	0.12	8.76
"sample"	0.06	4.38	0.18	13.14
"mean"	0.04	2.92	0.10	7.30
"mean.default"	0.03	2.19	0.06	4.38
"length"	0.03	2.19	0.03	2.19
"simplify2array"	0.01	0.73	0.02	1.46
"=="	0.01	0.73	0.01	0.73
"is.numeric"	0.01	0.73	0.01	0.73
"paste"	0.01	0.73	0.01	0.73
"unique"	0.01	0.73	0.01	0.73

\$by.total

	total.time	total.pct	self.time	self.pct
"block_exec"	1.37	100.00	0.00	0.00
"call_block"	1.37	100.00	0.00	0.00
"evaluate_call"	1.37	100.00	0.00	0.00
"evaluate::evaluate"	1.37	100.00	0.00	0.00
"evaluate"	1.37	100.00	0.00	0.00
"in_dir"	1.37	100.00	0.00	0.00
"knitr::knit"	1.37	100.00	0.00	0.00
"process_file"	1.37	100.00	0.00	0.00
"process_group.block"	1.37	100.00	0.00	0.00
"process_group"	1.37	100.00	0.00	0.00
"rmarkdown::render"	1.37	100.00	0.00	0.00
"withCallingHandlers"	1.37	100.00	0.00	0.00
"FUN"	1.36	99.27	1.05	76.64
"doTryCatch"	1.36	99.27	0.00	0.00
"eval"	1.36	99.27	0.00	0.00
"handle"	1.36	99.27	0.00	0.00
"lapply"	1.36	99.27	0.00	0.00
"sapply"	1.36	99.27	0.00	0.00
"timing_fn"	1.36	99.27	0.00	0.00
"try"	1.36	99.27	0.00	0.00
"tryCatch"	1.36	99.27	0.00	0.00
"tryCatchList"	1.36	99.27	0.00	0.00
"tryCatchOne"	1.36	99.27	0.00	0.00

"withVisible"	1.36	99.27	0.00	0.00
"sample"	0.18	13.14	0.06	4.38
"sample.int"	0.12	8.76	0.11	8.03
"mean"	0.10	7.30	0.04	2.92
"mean.default"	0.06	4.38	0.03	2.19
"length"	0.03	2.19	0.03	2.19
"simplify2array"	0.02	1.46	0.01	0.73
"=="	0.01	0.73	0.01	0.73
"is.numeric"	0.01	0.73	0.01	0.73
"paste"	0.01	0.73	0.01	0.73
"unique"	0.01	0.73	0.01	0.73
".make_numeric_version"	0.01	0.73	0.00	0.00
"getRversion"	0.01	0.73	0.00	0.00
"handle_output"	0.01	0.73	0.00	0.00
"package_version"	0.01	0.73	0.00	0.00
"plot_snapshot"	0.01	0.73	0.00	0.00
"R_system_version"	0.01	0.73	0.00	0.00
"recordPlot"	0.01	0.73	0.00	0.00
"w\$get_new"	0.01	0.73	0.00	0.00

```
$sample.interval
[1] 0.01
```

```
$sampling.time
[1] 1.37
```

Looping

That produced a minor, but useful improvement in speed. As we look through the rest of timings, we can see that there aren't a lot more savings to get. Most of the time is being taken by `sapply`, which we need in order to do the looping. We can use another member of the `apply` family to do grouped iterations: `tapply`:

```
df <- read.csv("ctd.csv", stringsAsFactors = FALSE)

# Start profiling
Rprof(interval = 0.01)

# Run code
df.1 <- df[df$depth == 1, ]
temp <- df.1$temp
station <- df.1$station
st.names <- unique(station)

boot.mean <- sapply(1:1000, function(i) {
  tapply(temp, station, function(st.temp) {
    st.temp <- sample(st.temp, length(st.temp), replace = TRUE)
    mean(st.temp)
  })
})
ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))

# Stop profiling
Rprof(NULL)
```

```
# Examine profile summary
summaryRprof()
```

```
$by.self
```

	self.time	self.pct	total.time	total.pct
"sample.int"	0.19	30.16	0.19	30.16
"sample"	0.07	11.11	0.26	41.27
"mean"	0.06	9.52	0.09	14.29
"FUN"	0.04	6.35	0.62	98.41
"lapply"	0.04	6.35	0.62	98.41
"as.character"	0.03	4.76	0.03	4.76
"mean.default"	0.03	4.76	0.03	4.76
"split.default"	0.03	4.76	0.03	4.76
"factor"	0.02	3.17	0.11	17.46
"sort.list"	0.02	3.17	0.04	6.35
"unique"	0.02	3.17	0.04	6.35
"unique.default"	0.02	3.17	0.02	3.17
"c"	0.01	1.59	0.01	1.59
"grepl"	0.01	1.59	0.01	1.59
"is.numeric"	0.01	1.59	0.01	1.59
"is.ordered"	0.01	1.59	0.01	1.59
"lazyLoadDBfetch"	0.01	1.59	0.01	1.59
"unlist"	0.01	1.59	0.01	1.59

```
$by.total
```

	total.time	total.pct	self.time	self.pct
"block_exec"	0.63	100.00	0.00	0.00
"call_block"	0.63	100.00	0.00	0.00
"evaluate_call"	0.63	100.00	0.00	0.00
"evaluate::evaluate"	0.63	100.00	0.00	0.00
"evaluate"	0.63	100.00	0.00	0.00
"in_dir"	0.63	100.00	0.00	0.00
"knitr::knit"	0.63	100.00	0.00	0.00
"process_file"	0.63	100.00	0.00	0.00
"process_group.block"	0.63	100.00	0.00	0.00
"process_group"	0.63	100.00	0.00	0.00
"rmarkdown::render"	0.63	100.00	0.00	0.00
"withCallingHandlers"	0.63	100.00	0.00	0.00
"FUN"	0.62	98.41	0.04	6.35
"lapply"	0.62	98.41	0.04	6.35
"doTryCatch"	0.62	98.41	0.00	0.00
"eval"	0.62	98.41	0.00	0.00
"handle"	0.62	98.41	0.00	0.00
"sapply"	0.62	98.41	0.00	0.00
"timing_fn"	0.62	98.41	0.00	0.00
"try"	0.62	98.41	0.00	0.00
"tryCatch"	0.62	98.41	0.00	0.00
"tryCatchList"	0.62	98.41	0.00	0.00
"tryCatchOne"	0.62	98.41	0.00	0.00
"withVisible"	0.62	98.41	0.00	0.00
"tapply"	0.61	96.83	0.00	0.00
"sample"	0.26	41.27	0.07	11.11
"sample.int"	0.19	30.16	0.19	30.16
"factor"	0.11	17.46	0.02	3.17

"mean"	0.09	14.29	0.06	9.52
"sort.list"	0.04	6.35	0.02	3.17
"unique"	0.04	6.35	0.02	3.17
"as.character"	0.03	4.76	0.03	4.76
"mean.default"	0.03	4.76	0.03	4.76
"split.default"	0.03	4.76	0.03	4.76
"split"	0.03	4.76	0.00	0.00
"unique.default"	0.02	3.17	0.02	3.17
"c"	0.01	1.59	0.01	1.59
"grepl"	0.01	1.59	0.01	1.59
"is.numeric"	0.01	1.59	0.01	1.59
"is.ordered"	0.01	1.59	0.01	1.59
"lazyLoadDBfetch"	0.01	1.59	0.01	1.59
"unlist"	0.01	1.59	0.01	1.59
".make_numeric_version"	0.01	1.59	0.00	0.00
"array"	0.01	1.59	0.00	0.00
"force"	0.01	1.59	0.00	0.00
"getRversion"	0.01	1.59	0.00	0.00
"handle_output"	0.01	1.59	0.00	0.00
"match.arg"	0.01	1.59	0.00	0.00
"package_version"	0.01	1.59	0.00	0.00
"plot_snapshot"	0.01	1.59	0.00	0.00
"R_system_version"	0.01	1.59	0.00	0.00
"recordPlot"	0.01	1.59	0.00	0.00
"w\$get_new"	0.01	1.59	0.00	0.00

```
$sample.interval
[1] 0.01
```

```
$sampling.time
[1] 0.63
```

Although we're still spending time in the `tapply` and `sapply`, we've cut down the total time considerably. One thing we can check is if there is an effect of the order of the loops. Currently, we are calculating the mean for all stations for each replicate. Lets switch the order so that we are calculating the means of all replicates for each station:

```
df <- read.csv("ctd.csv", stringsAsFactors = FALSE)

# Start profiling
Rprof(interval = 0.01)

# Run code
df.1 <- df[df$depth == 1, ]
temp <- df.1$temp
station <- df.1$station
st.names <- unique(station)

boot.mean <- tapply(temp, station, function(x) {
  sapply(1:1000, function(i) {
    st.temp <- sample(x, length(x), replace = TRUE)
    mean(st.temp)
  })
})
boot.mean <- do.call(rbind, boot.mean)
```



```
ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))
```

```
# Stop profiling
```

```
Rprof(NULL)
```

```
# Examine profile summary
```

```
summaryRprof()
```

```
$by.self
```

	self.time	self.pct	total.time	total.pct
"sample.int"	0.14	42.42	0.14	42.42
"sample"	0.07	21.21	0.21	63.64
"FUN"	0.04	12.12	0.32	96.97
"mean"	0.03	9.09	0.05	15.15
"lapply"	0.02	6.06	0.32	96.97
"mean.default"	0.01	3.03	0.02	6.06
"is.numeric"	0.01	3.03	0.01	3.03
"setHook"	0.01	3.03	0.01	3.03

```
$by.total
```

	total.time	total.pct	self.time	self.pct
"block_exec"	0.33	100.00	0.00	0.00
"call_block"	0.33	100.00	0.00	0.00
"evaluate_call"	0.33	100.00	0.00	0.00
"evaluate::evaluate"	0.33	100.00	0.00	0.00
"evaluate"	0.33	100.00	0.00	0.00
"in_dir"	0.33	100.00	0.00	0.00
"knitr::knit"	0.33	100.00	0.00	0.00
"process_file"	0.33	100.00	0.00	0.00
"process_group.block"	0.33	100.00	0.00	0.00
"process_group"	0.33	100.00	0.00	0.00
"rmarkdown::render"	0.33	100.00	0.00	0.00
"withCallingHandlers"	0.33	100.00	0.00	0.00
"FUN"	0.32	96.97	0.04	12.12
"lapply"	0.32	96.97	0.02	6.06
"doTryCatch"	0.32	96.97	0.00	0.00
"eval"	0.32	96.97	0.00	0.00
"handle"	0.32	96.97	0.00	0.00
"sapply"	0.32	96.97	0.00	0.00
"tapply"	0.32	96.97	0.00	0.00
"timing_fn"	0.32	96.97	0.00	0.00
"try"	0.32	96.97	0.00	0.00
"tryCatch"	0.32	96.97	0.00	0.00
"tryCatchList"	0.32	96.97	0.00	0.00
"tryCatchOne"	0.32	96.97	0.00	0.00
"withVisible"	0.32	96.97	0.00	0.00
"sample"	0.21	63.64	0.07	21.21
"sample.int"	0.14	42.42	0.14	42.42
"mean"	0.05	15.15	0.03	9.09
"mean.default"	0.02	6.06	0.01	3.03
"is.numeric"	0.01	3.03	0.01	3.03
"setHook"	0.01	3.03	0.01	3.03
"set_hooks"	0.01	3.03	0.00	0.00

```
$sample.interval  
[1] 0.01
```

```
$sampling.time  
[1] 0.33
```

This is slightly faster because we are not doing the `tapply` 1000 times, which takes some time. It doesn't look like we can do much more. However, there is a more efficient way of looping, although it is not necessarily as compact. Instead of using the interior `sapply` constructs, we can pre-allocate a result vector, and use a `for` loop to fill it:

```
df <- read.csv("ctd.csv", stringsAsFactors = FALSE)  
  
# Start profiling  
Rprof(interval = 0.01)  
  
# Run code  
df.1 <- df[df$depth == 1, ]  
temp <- df.1$temp  
station <- df.1$station  
st.names <- unique(station)  
  
# an empty result vector  
boot.vec <- vector(length = 1000)  
  
boot.mean <- tapply(temp, station, function(x) {  
  for(i in 1:length(boot.vec)) {  
    st.temp <- sample(x, length(x), replace = TRUE)  
    boot.vec[i] <- mean(st.temp)  
  }  
  boot.vec  
)  
boot.mean <- do.call(rbind, boot.mean)  
ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))  
  
# Stop profiling  
Rprof(NULL)  
  
# Examine profile summary  
summaryRprof()
```

```
$by.self  
      self.time self.pct total.time total.pct  
"sample"      0.15  45.45      0.23  69.70  
"sample.int"  0.07  21.21      0.08  24.24  
"mean"        0.06  18.18      0.08  24.24  
"FUN"         0.01   3.03      0.33 100.00  
"mean.default" 0.01   3.03      0.02   6.06  
"("          0.01   3.03      0.01   3.03  
"findCenvVar" 0.01   3.03      0.01   3.03  
"is.numeric"  0.01   3.03      0.01   3.03  
  
$by.total  
      total.time total.pct self.time self.pct  
"FUN"          0.33  100.00      0.01   3.03
```

"block_exec"	0.33	100.00	0.00	0.00
"call_block"	0.33	100.00	0.00	0.00
"doTryCatch"	0.33	100.00	0.00	0.00
"eval"	0.33	100.00	0.00	0.00
"evaluate_call"	0.33	100.00	0.00	0.00
"evaluate::evaluate"	0.33	100.00	0.00	0.00
"evaluate"	0.33	100.00	0.00	0.00
"handle"	0.33	100.00	0.00	0.00
"in_dir"	0.33	100.00	0.00	0.00
"knitr::knit"	0.33	100.00	0.00	0.00
"lapply"	0.33	100.00	0.00	0.00
"process_file"	0.33	100.00	0.00	0.00
"process_group.block"	0.33	100.00	0.00	0.00
"process_group"	0.33	100.00	0.00	0.00
"rmarkdown::render"	0.33	100.00	0.00	0.00
"tapply"	0.33	100.00	0.00	0.00
"timing_fn"	0.33	100.00	0.00	0.00
"try"	0.33	100.00	0.00	0.00
"tryCatch"	0.33	100.00	0.00	0.00
"tryCatchList"	0.33	100.00	0.00	0.00
"tryCatchOne"	0.33	100.00	0.00	0.00
"withCallingHandlers"	0.33	100.00	0.00	0.00
"withVisible"	0.33	100.00	0.00	0.00
"sample"	0.23	69.70	0.15	45.45
"sample.int"	0.08	24.24	0.07	21.21
"mean"	0.08	24.24	0.06	18.18
"mean.default"	0.02	6.06	0.01	3.03
"("	0.01	3.03	0.01	3.03
"findCenvVar"	0.01	3.03	0.01	3.03
"is.numeric"	0.01	3.03	0.01	3.03
"checkSkipLoopCntxt"	0.01	3.03	0.00	0.00
"cmp"	0.01	3.03	0.00	0.00
"cmpCall"	0.01	3.03	0.00	0.00
"cmpfun"	0.01	3.03	0.00	0.00
"compiler::tryCmpfun"	0.01	3.03	0.00	0.00
"genCode"	0.01	3.03	0.00	0.00
"getInlineInfo"	0.01	3.03	0.00	0.00
"h"	0.01	3.03	0.00	0.00
"isBaseVar"	0.01	3.03	0.00	0.00
"isLoopTopFun"	0.01	3.03	0.00	0.00
"tryInline"	0.01	3.03	0.00	0.00

```
$sample.interval
[1] 0.01
```

```
$sampling.time
[1] 0.33
```

Well, that was faster overall. We can now extend the concept and try the same thing for the tapply loop. This time, we have to create a matrix to hold the results.

```
df <- read.csv("ctd.csv", stringsAsFactors = FALSE)

# Start profiling
Rprof(interval = 0.01)
```

```

# Run code
df.1 <- df[df$depth == 1, ]
temp <- df.1$temp
station <- df.1$station
st.names <- unique(station)

# an empty result vector
boot.mean <- matrix(nrow = length(st.names), ncol = 1000)
rownames(boot.mean) <- st.names

for(st in st.names) {
  x <- temp[station == st]
  num.temp <- length(x)
  for(i in 1:ncol(boot.mean)) {
    st.temp <- sample(x, num.temp, replace = TRUE)
    boot.mean[st, i] <- mean(st.temp)
  }
}
ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))

# Stop profiling
Rprof(NULL)

# Examine profile summary
summaryRprof()

```

```

$by.self
      self.time self.pct total.time total.pct
"sample"      0.09  23.68      0.19   50.00
"sample.int"   0.08  21.05      0.08   21.05
"eval"         0.05  13.16      0.38  100.00
"mean"         0.05  13.16      0.12   31.58
"mean.default" 0.05  13.16      0.07   18.42
"is.numeric"   0.02   5.26      0.02   5.26
"length"       0.02   5.26      0.02   5.26
"formatC"      0.01   2.63      0.01   2.63
"h"            0.01   2.63      0.01   2.63

$by.total
      total.time total.pct self.time self.pct
"eval"          0.38  100.00      0.05  13.16
"block_exec"    0.38  100.00      0.00   0.00
"call_block"    0.38  100.00      0.00   0.00
"doTryCatch"    0.38  100.00      0.00   0.00
"evaluate_call" 0.38  100.00      0.00   0.00
"evaluate::evaluate" 0.38  100.00      0.00   0.00
"evaluate"      0.38  100.00      0.00   0.00
"handle"        0.38  100.00      0.00   0.00
"in_dir"        0.38  100.00      0.00   0.00
"knitr::knit"   0.38  100.00      0.00   0.00
"process_file"  0.38  100.00      0.00   0.00
"process_group.block" 0.38  100.00      0.00   0.00
"process_group" 0.38  100.00      0.00   0.00

```

"rmarkdown::render"	0.38	100.00	0.00	0.00
"timing_fn"	0.38	100.00	0.00	0.00
"try"	0.38	100.00	0.00	0.00
"tryCatch"	0.38	100.00	0.00	0.00
"tryCatchList"	0.38	100.00	0.00	0.00
"tryCatchOne"	0.38	100.00	0.00	0.00
"withCallingHandlers"	0.38	100.00	0.00	0.00
"withVisible"	0.38	100.00	0.00	0.00
"sample"	0.19	50.00	0.09	23.68
"mean"	0.12	31.58	0.05	13.16
"sample.int"	0.08	21.05	0.08	21.05
"mean.default"	0.07	18.42	0.05	13.16
"is.numeric"	0.02	5.26	0.02	5.26
"length"	0.02	5.26	0.02	5.26
"formatC"	0.01	2.63	0.01	2.63
"h"	0.01	2.63	0.01	2.63
"apply"	0.01	2.63	0.00	0.00
"cmp"	0.01	2.63	0.00	0.00
"cmpCall"	0.01	2.63	0.00	0.00
"compile"	0.01	2.63	0.00	0.00
"compiler:::tryCompile"	0.01	2.63	0.00	0.00
"format_perc"	0.01	2.63	0.00	0.00
"FUN"	0.01	2.63	0.00	0.00
"genCode"	0.01	2.63	0.00	0.00
"paste0"	0.01	2.63	0.00	0.00
"quantile.default"	0.01	2.63	0.00	0.00
"tryInline"	0.01	2.63	0.00	0.00

```
$sample.interval
[1] 0.01
```

```
$sampling.time
[1] 0.38
```

That pre-allocation speeds us up even more. The only other things we can do is use the function `sample.int` directly rather than through `sample`, and use the internal function `.Internal(mean())` rather than the generic `mean`:

```
df <- read.csv("ctd.csv", stringsAsFactors = FALSE)

# Start profiling
Rprof(interval = 0.01)

# Run code
df.1 <- df[df$depth == 1, ]
temp <- df.1$temp
station <- df.1$station
st.names <- unique(station)

# an empty result vector
boot.mean <- matrix(nrow = length(st.names), ncol = 1000)
rownames(boot.mean) <- st.names

for(st in st.names) {
  x <- temp[station == st]
```

```

num.temp <- length(x)
for(i in 1:ncol(boot.mean)) {
  j <- sample.int(1:num.temp, num.temp, replace = TRUE)
  boot.mean[st, i] <- .Internal(mean(x[j]))
}
}
ci <- apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))

# Stop profiling
Rprof(NULL)

# Examine profile summary
summaryRprof()

$by.self
              self.time self.pct total.time total.pct
"sample.int"      0.13   61.90      0.13   61.90
"eval"            0.07   33.33      0.20   95.24
".make_numeric_version" 0.01    4.76      0.01    4.76

$by.total
              total.time total.pct self.time self.pct
"block_exec"      0.21   100.00      0.00    0.00
"call_block"      0.21   100.00      0.00    0.00
"evaluate_call"   0.21   100.00      0.00    0.00
"evaluate::evaluate" 0.21   100.00      0.00    0.00
"evaluate"        0.21   100.00      0.00    0.00
"in_dir"          0.21   100.00      0.00    0.00
"knitr::knit"     0.21   100.00      0.00    0.00
"process_file"    0.21   100.00      0.00    0.00
"process_group.block" 0.21   100.00      0.00    0.00
"process_group"   0.21   100.00      0.00    0.00
"rmarkdown::render" 0.21   100.00      0.00    0.00
"withCallingHandlers" 0.21   100.00      0.00    0.00
"eval"            0.20   95.24      0.07   33.33
"doTryCatch"      0.20   95.24      0.00    0.00
"handle"          0.20   95.24      0.00    0.00
"timing_fn"        0.20   95.24      0.00    0.00
"try"             0.20   95.24      0.00    0.00
"tryCatch"        0.20   95.24      0.00    0.00
"tryCatchList"    0.20   95.24      0.00    0.00
"tryCatchOne"     0.20   95.24      0.00    0.00
"withVisible"     0.20   95.24      0.00    0.00
"sample.int"      0.13   61.90      0.13   61.90
".make_numeric_version" 0.01    4.76      0.01    4.76
"getRversion"     0.01    4.76      0.00    0.00
"handle_output"   0.01    4.76      0.00    0.00
"package_version" 0.01    4.76      0.00    0.00
"plot_snapshot"   0.01    4.76      0.00    0.00
"R_system_version" 0.01    4.76      0.00    0.00
"recordPlot"      0.01    4.76      0.00    0.00
"w$get_new"        0.01    4.76      0.00    0.00

$sample.interval

```

```
[1] 0.01
```

```
$sampling.time
```

```
[1] 0.21
```

Using `.Internal` functions can be tricky because the code base of internal functions can change over time. Also, CRAN will not permit packages using `.Internal` to be submitted.

Benchmarking

We've seen incremental achievements in our code, but it would be good to know how much better one version of code is than another. There are a couple of ways to do this. The first is to use `system.time` to record the CPU time required for a set of expressions to execute. For our comparison, let's make three functions that represent our code at different stages and see how long each actually takes:

```
# the first one
bootMean.1 <- function(fname, nrep) {
  boot.mean <- sapply(1:nrep, function(i) {
    df <- read.csv(fname, stringsAsFactors = FALSE)
    sapply(unique(df$station), function(st) {
      i <- which(df$station == st & df$depth == 1)
      i <- sample(i, length(i), replace = TRUE)
      mean(df$temp[i])
    })
  })
  apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))
}

# halfway through optimizing, using two sapply functions
bootMean.2 <- function(fname, nrep) {
  df <- read.csv(fname, stringsAsFactors = FALSE)
  boot.mean <- sapply(1:nrep, function(i) {
    sapply(unique(df$station), function(st) {
      i <- which(df$station == st & df$depth == 1)
      i <- sample(i, length(i), replace = TRUE)
      mean(df$temp[i])
    })
  })
  apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))
}

# fully optimized, using for-loops
bootMean.3 <- function(fname, nrep) {
  df <- read.csv(fname, stringsAsFactors = FALSE)
  df <- df[df$depth == 1, ]
  temp <- df$temp
  station <- df$station
  st.names <- unique(station)

  boot.mean <- matrix(nrow = length(st.names), ncol = nrep)
  rownames(boot.mean) <- st.names

  for(st in st.names) {
```

```

x <- temp[station == st]
num.temp <- length(x)
for(i in 1:ncol(boot.mean)) {
  j <- sample.int(1:num.temp, num.temp, replace = TRUE)
  boot.mean[st, i] <- .Internal(mean(x[j]))
}
}
apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))
}

# CPU time for each function for 10 replicates
system.time(bootMean.1("ctd.csv", 10))

```

```

user  system elapsed
5.940  0.212   6.163

```

```
system.time(bootMean.2("ctd.csv", 10))
```

```

user  system elapsed
0.991  0.122   1.114

```

```
system.time(bootMean.3("ctd.csv", 10))
```

```

user  system elapsed
0.619  0.009   0.629

```

It looks like we were able to make it about 10 times faster. Another useful way to do this is to use the `microbenchmark` package, which will execute each expression a number of times and give a distribution of the run times.

```

library(microbenchmark)
microbenchmark(
  bootMean.1 = bootMean.1("ctd.csv", 10),
  bootMean.2 = bootMean.2("ctd.csv", 10),
  bootMean.3 = bootMean.3("ctd.csv", 10),
  times = 10
)

```

Unit: milliseconds

	expr	min	lq	mean	median	uq	max
bootMean.1		6130.5694	6258.5644	6436.3778	6492.3210	6567.8607	6680.0761
bootMean.2		1066.9551	1083.7965	1125.0335	1143.0086	1146.1704	1198.0278
bootMean.3		552.4734	579.7782	606.3579	617.7394	630.6384	645.8263
neval							
	10						
	10						
	10						

Parallel computing

Once we have code that is optimized as much as possible in R, if we need it to execute faster, we can explore distributing the processing among several CPUs if we have access to a computer that has multiple cores. This is only good if we have processes that are independent, such that if one is running on one core, it does not need to transfer information to a process on another core. This is good for simulations, permutations, bootstrapping, and other similar loop-based operations. There is a bit of overhead in setting up the monitoring

of the parallel processes, which takes some time, so make sure that this will be made up for in the time it takes to run processes in parallel rather than serially.

The easiest way to get into parallel programming is to use the `parallel` package. We should first find out how many cores we have:

```
library(parallel)
detectCores()
```

```
[1] 8
```

The way my MacBook is configured, it will report `{r}` `detectCores()` cores. However, that is not totally true. It actually has half that number as each physical CPU is virtually doubled. While I can allocate that many cores for independent work, it wouldn't leave any space on my CPU for overhead processes. Therefore, I like to use no more than `{r}` $(\text{detectCores()} / 2) - 1$ cores when I set up intensive processes.

Below, we'll use the `tapply` version of our example and structure the code to spread the replicates for a station among 2 cores.

```
bootMean.tapply <- function(fname, nrep) {
  df <- read.csv(fname, stringsAsFactors = FALSE)
  df <- df[df$depth == 1, ]
  temp <- df$temp
  station <- df$station
  st.names <- unique(station)

  boot.mean <- tapply(temp, station, function(x) {
    sapply(1:nrep, function(i) {
      st.temp <- sample(x, length(x), replace = TRUE)
      mean(st.temp)
    })
  })
  boot.mean <- do.call(rbind, boot.mean)
  apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))
}
```

```
bootMean.par <- function(fname, nrep) {
  df <- read.csv(fname, stringsAsFactors = FALSE)
  df <- df[df$depth == 1, ]
  temp <- df$temp
  station <- df$station
  st.names <- unique(station)

  # make 2 clusters
  cl <- makeCluster(2)

  boot.mean <- tapply(temp, station, function(x) {
    clusterExport(cl, "x")
    parSapply(cl, 1:nrep, function(i) {
      st.temp <- sample(x, length(x), replace = TRUE)
      mean(st.temp)
    })
  })
  stopCluster(cl)
  boot.mean <- do.call(rbind, boot.mean)
  apply(boot.mean, 1, quantile, probs = c(0.025, 0.975))
}
```

Here's the result for 100 replicates:

```
system.time(bootMean.tapply("ctd.csv", 100))
```

```
   user  system elapsed  
0.667   0.016   0.688
```

```
system.time(bootMean.par("ctd.csv", 100))
```

```
   user  system elapsed  
0.829   0.063   1.363
```

Here's the result for 1000 replicates:

```
system.time(bootMean.tapply("ctd.csv", 1000))
```

```
   user  system elapsed  
0.963   0.024   0.989
```

```
system.time(bootMean.par("ctd.csv", 1000))
```

```
   user  system elapsed  
0.781   0.051   1.482
```

Here's the result for 10,000 replicates:

```
system.time(bootMean.tapply("ctd.csv", 10000))
```

```
   user  system elapsed  
4.928   0.148   5.099
```

```
system.time(bootMean.par("ctd.csv", 10000))
```

```
   user  system elapsed  
1.043   0.089   3.517
```

You can see that at 100 replicates, the parallel version takes more time than the non-parallel version. However, as the number of replicates increases, the parallel version gets much faster. You can see that it would be very useful to parallelize code if each replicate was even moderately time-intensive.