

Lecture 24: Profiling and microbenchmarking

An order of operations for programming

1. Make it run
2. Make it right
3. Make it fast

When speed matters

- You are working with very large data
- You are running a process (a simulation, a data analysis, etc.) many times
- A piece of code will be called many times (e.g., choosing a split in a decision tree)

Goals

- Learn how to identify bottlenecks in code
- Learn approaches for more efficient code in R
- Time permitting: learn how to use C++ to make code faster

Example: timing code

Suppose we want to compute the mean of each column of a data frame:

```
1  n <- 100000
2  cols <- 150
3  data_mat <- matrix(rnorm(n * cols, mean = 5), ncol = cols)
4  data <- as.data.frame(data_mat)
5
6  means <- rep(NA, cols)
7  for(i in 1:cols){
8    means[i] <- mean(data[,i])
9  }
```

Example: timing code

Suppose we want to compute the mean of each column of a data frame:

```
1  n <- 100000
2  cols <- 150
3  data_mat <- matrix(rnorm(n * cols, mean = 5), ncol = cols)
4  data <- as.data.frame(data_mat)
5
6  system.time({
7    means <- rep(NA, cols)
8    for(i in 1:cols){
9      means[i] <- mean(data[,i])
10   }
11 })
```

← timing code

user	system	elapsed	(in seconds)
1.923	0.014	1.939	

Alternatives

```
1 means <- rep(NA, cols)
2 for(i in 1:cols){
3   means[i] <- mean(data[,i])
4 }
```

What are the alternatives to this for-loop approach?

- Is there a built-in function?

colMeans : function for computing
column means of matrices
& data frames

- Is there another way to iterate?

apply : apply a function to the
margins of a matrix or data
frame

apply(data, 2, mean)

Alternatives

```
1 # Option 1: for loop
2 for_loop_means <- function(data){
3   cols <- ncol(data)
4   means <- rep(NA, cols)
5   for(i in 1:cols){
6     means[i] <- mean(data[,i])
7   }
8   return(means)
9 }
10 means <- for_loop_means(data)
11
12 # Option 2: apply
13 means <- apply(data, 2, mean)
14
15 # Option 3: colMeans
16 means <- colMeans(data)
```


Comparing performance

Microbenchmarking: Evaluating the performance of a small piece of code

useful for benchmarking

```
1 bench::mark(  
2   means <- for_loop_means(data),  
3   means <- apply(data, 2, mean),  
4   means <- colMeans(data),  
5   check = F  
6 )
```

A tibble: 3 × 6

expression	min	median	`itr/sec`	mem_alloc
`gc/sec`				
<bch:expr>	<bch:tm>	<bch:tm>	<dbl>	<bch:byt>
<dbl>				
1 means <- for_loop_means(data)	1.92s	1.92s	0.522	1.85KB
0				
2 means <- apply(data, 2, mean)	2.01s	2.01s	0.497	400.57MB
0.993				
3 means <- colMeans(data)	450.29ms	462.51ms	2.16	114.45MB
1.08				

least memory used

most memory used

time it took for the code to run once

about 4 times faster

similar performance

Profiling

```
1 library(profvis)
2 profvis({
3   means <- for_loop_means(data)
4   means <- apply(data, 2, mean)
5   means <- colMeans(data)
6 })
```

```
1 profvis({
2   means <- for_loop_means(data)
3   means <- apply(data, 2, mean)
4   means <- colMeans(data)
5 })
```

memory used time (ms)

	1790
400.6	1840
114.5	480

converting data
frame into a matrix

as.matrix

aperm

mean.default

as.matrix

mean.default
for_loop_means

mean.default

apply

colMeans

mean.default
(this takes
almost all
the time
in the
for loop)

for loop

apply

colmeans

(length reflects time spent)

Space for efficiency increases?

```
1 colMeans
```

```
function (x, na.rm = FALSE, dims = 1L)
{
  if (is.data.frame(x))
    x <- as.matrix(x)
  if (!is.array(x) || length(dn <- dim(x)) < 2L)
    stop("'x' must be an array of at least two dimensions")
  if (dims < 1L || dims > length(dn) - 1L)
    stop("invalid 'dims'")
  n <- prod(dn[id <- seq_len(dims)])
  dn <- dn[-id]
  z <- if (is.complex(x))
    .Internal(colMeans(Re(x), n, prod(dn), na.rm)) + (0+1i) *
      .Internal(colMeans(Im(x), n, prod(dn), na.rm))
  else .Internal(colMeans(x, n, prod(dn), na.rm))
}
```

checking
inputs,
doing
transformations

Increase efficiency by avoiding extraneous steps

```
1 n <- 100000
2 cols <- 150
3 data_mat <- matrix(rnorm(n * cols, mean = 5), ncol = cols)
4 data <- as.data.frame(data_mat)
5
6 bench::mark(
7   means <- colMeans(data_mat),
8   means <- colMeans(data),
9   check = F
10 )
```

A tibble: 2 × 6

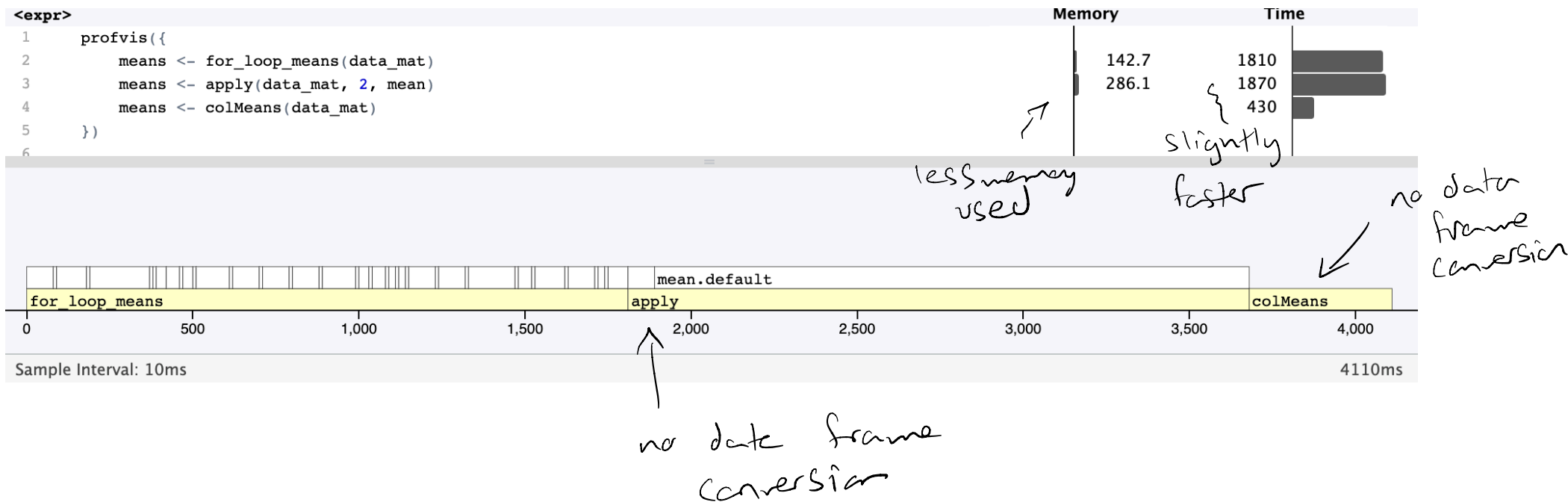
expression	min	median	`itr/sec`	mem_alloc
<code>`gc/sec`</code>				
<code><bch:expr></code>	<code><bch:tm></code>	<code><bch:tm></code>	<code><dbl></code>	<code><bch:byt></code>
<code><dbl></code>				
1 means <- colMeans(data_mat)	434ms	436ms	2.29	1.22KB
0				
2 means <- colMeans(data)	450ms	450ms	2.22	114.45MB
2.22				

less
memory
used

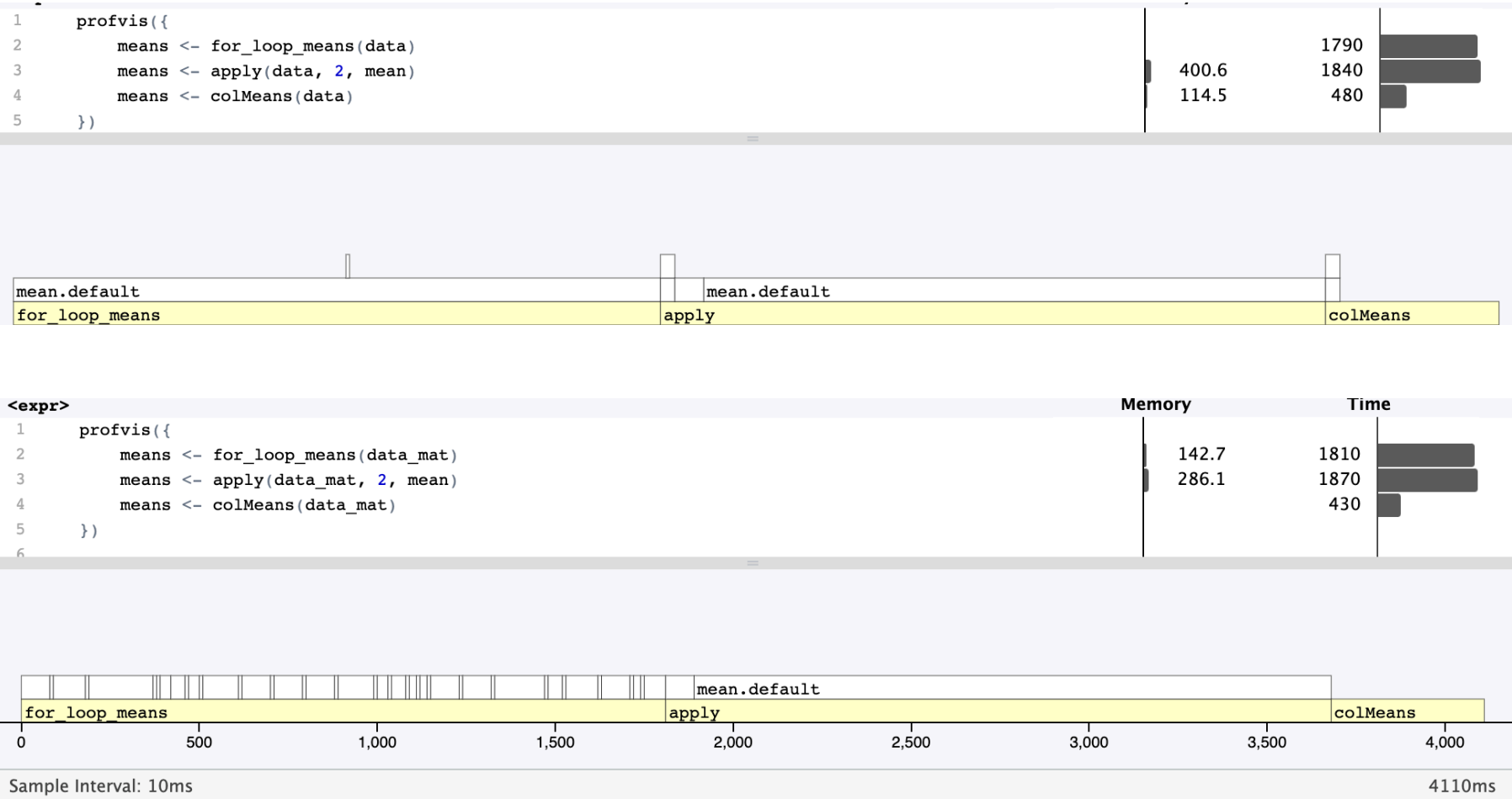
slightly
faster

Profiling

```
1 profvis({  
2   means <- for_loop_means(data_mat)  
3   means <- apply(data_mat, 2, mean)  
4   means <- colMeans(data_mat)  
5 })
```



Profiling



Class activity

https://sta279-s24.github.io/class_activities/ca_lecture_24.html

