

Lecture 25: Making code more efficient

Approaches to faster code

- Do as little as possible
- Vectorise
- Avoid copies

Do as little as possible

(remove extraneous step
if you're sure they are
extraneous)

```
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><bch:expr></code>	<code><bch:tm></code>	<code><bch:tm></code>	<code><dbl></code>	<code><bch:byt></code>
1 means <- colMeans(data_mat)	437ms	437ms	2.29	25.4KB
2 means <- colMeans(data)	456ms	458ms	2.18	114.5MB

Avoid copies

The code below samples 100 observations from a $N(0, 1)$ distribution:

```
1 x <- c()
2 for(i in 1:100){
3   x <- c(x, rnorm(1))
4 }
```

How could I make this code more efficient?

- Growing vector is slow
(every time I add a new entry I have to
create a new object $\hat{=}$ allocate memory)

Alternative: create a vector of the right length
then fill it in

Avoid copies

```
1 loop_1 <- function(n){
2   x <- c()
3   for(i in 1:n){
4     x <- c(x, rnorm(1))
5   }
6   return(x)
7 }
8
9 loop_2 <- function(n){
10  x <- rep(NA, n)
11  for(i in 1:n){
12    x[i] <- rnorm(1)
13  }
14  return(x)
15 }
```

Avoid copies

```
1 bench::mark(  
2   loop_1(100),  
3   loop_2(100),  
4   check = F  
5 )
```

A tibble: 2 × 6

	expression	min	median	`itr/sec`	mem_alloc	`gc/sec`
	<bch:expr>	<bch:tm>	<bch:tm>	<dbl>	<bch:byt>	<dbl>
1	loop_1(100)	143µs	161µs	5552.	318KB	11.6
2	loop_2(100)	114µs	118µs	7529.	272KB	11.5

↑
slightly faster

↑
less memory

Avoid copies

```
1 bench::mark(  
2   loop_1(10000),  
3   loop_2(10000),  
4   check = F  
5 )
```

```
# A tibble: 2 × 6
```

expression	min	median	`itr/sec`	mem_alloc	`gc/sec`
<bch:expr>	<bch:tm>	<bch:tm>	<dbl>	<bch:byt>	<dbl>
1 loop_1(10000)	85.7ms	104.8ms	9.22	406.3MB	20.0
2 loop_2(10000)	11.8ms	12.2ms	70.6	24.5MB	9.80

↑
much faster

↑
much less memory

Vectorise

The code below samples 100 observations from a $N(0, 1)$ distribution:

```
1 x <- rep(NA, 100)
2 for(i in 1:100){
3   x[i] <- rnorm(1)
4 }
```

How could I make this code more efficient?

rnorm(100)

Vectorise

```
1 for_loop_sample <- function(n){
2   x <- rep(NA, n)
3   for(i in 1:n){
4     x[i] <- rnorm(1)
5   }
6 }
7
8 bench::mark(
9   x <- for_loop_sample(100),
10  x <- rnorm(100),
11  check=F
12 )
```

A tibble: 2 × 6

expression	min	median	`itr/sec`	mem_alloc
`gc/sec`				
<bch:expr>	<bch:tm>	<bch:tm>	<dbl>	<bch:byt>
<dbl>				
1 x <- for_loop_sample(100)	113.92µs	119.04µs	8069.	271.04KB
11.2				
2 x <- rnorm(100)	5.42µs	6.33µs	151517.	3.32KB
15.2				

↑
about 20x faster

Other options

- Different data structures / algorithms
- Parallelization
- Rewrite code in C++

Class activity

<https://sta279->

[f23.github.io/class_activities/ca_lecture_25.html](https://sta279-f23.github.io/class_activities/ca_lecture_25.html)

