

Class 6: R Functions

Kaliyah (A18125684)

Table of contents

Background	1
Our first function	1
A second function	2
A Protein generating function	5

Background

All functions in R have at least 3 things: - A **name** that we use to call the function. - One or more **arguments** - The **body** the lines of R code that do the work

Our first function

Let's write a silly wee function called `add()` to add some numbers (the input arguments)

```
add <- function(x,y) {  
  x+y  
}
```

Now we can use this function

```
add(100,1)
```

```
[1] 101
```

Q. What if I give a multiple element vector to x and y?

```
add(x=c(100,1), y=c(100,1))
```

```
[1] 200 2
```

What if I give three inputs to the function?

```
#add(x=c(100,1), y= 1 z=1)
```

Q. What if I give only one input to the add function

```
addnew <- function(x,y=1) {  
  x+y  
}
```

```
addnew(x=100)
```

```
[1] 101
```

```
addnew(c(100,1),100)
```

```
[1] 200 101
```

If we write our function with input arguments having no default value then the user will be required to set them when they use the function. We can give our input arguments “default” values by setting them equal to some sensible value. e.g x=1, y=1

A second function

Let's try something more interesting: Make a sequence generating tool...

The `sample()` function can be a useful starting point here:

```
sample (1:10, size=4)
```

```
[1] 8 3 1 4
```

Q. Generate 9 random numbers from the input vector x=1:10?

```
sample(1:10, size=9)
```

```
[1] 5 9 4 3 6 8 1 10 7
```

Q. Generate 12 random numbers from the input vector x=1:10?

```
sample(1:10, size =12, replace=TRUE)
```

```
[1] 7 1 4 5 1 4 4 4 4 2 4 2
```

Q. Write code for the `sample()` function that generates nucleotide sequences of length 6?

```
sample(x=c("A","C","G","T"), size=6, replace=TRUE )
```

```
[1] "T" "T" "C" "A" "C" "T"
```

Q. Write a first function `generate_dna()` that returns a *user specified length* DNA sequence:

```
generate_dna <- function(len=6) {  
  sample(x=c("A","C","G","T"), size = len, replace=TRUE )}
```

```
generate_dna(10)
```

```
[1] "T" "A" "G" "G" "T" "C" "T" "A" "C" "G"
```

```
generate_dna()
```

```
[1] "C" "G" "G" "C" "C" "C"
```

```
generate_dna(100)
```

```
[1] "G" "T" "G" "G" "A" "C" "C" "C" "T" "G" "T" "G" "T" "C" "C" "A" "A" "C"  
[19] "T" "G" "A" "G" "A" "C" "A" "A" "G" "C" "T" "C" "G" "A" "T" "A" "A" "A"  
[37] "C" "A" "C" "G" "A" "C" "G" "A" "G" "C" "G" "T" "T" "C" "G" "C" "G" "C" "G"  
[55] "G" "G" "C" "G" "A" "T" "G" "T" "C" "A" "C" "C" "A" "C" "A" "A" "A" "A"  
[73] "A" "G" "A" "G" "G" "C" "A" "A" "A" "G" "T" "T" "G" "G" "C" "A" "G" "T"  
[91] "G" "G" "A" "A" "G" "G" "T" "C" "C" "T"
```

Key Points Every function in R looks fundamentally the same in terms of its structure. Basically 3 things: name, input, and body

```
name <- function (input) {  
  body  
}
```

Functions can have multiple inputs. These can be **required** arguments or **optional** arguments. With optional arguments having a set default value.

Q. Modify and improve our `generate_dna()` function to return it's generated sequence in a more standard format like "AGTAGTA" rather than the vector "A", "C", "G", "A".

```
generate_dna <- function(len=6, fasta=TRUE) {  
  ans <- sample(x=c("A","C","G","T"),  
                size = len, replace=TRUE)  
  if (fasta) {  
    cat(" Single-element vector output")  
    ans <- paste(ans, collapse="")  
  } else {  
    cat("Multi-element vector output")  
  }  
  return(ans)  
}  
  
generate_dna(fasta=TRUE)
```

Single-element vector output

```
[1] "CGGTGA"
```

The `paste()` function - it's job is to join up or stick together (a.k.a paste) input strings together

```
paste(c("alice","barry"), "loves R", sep= " ")
```

```
[1] "alice loves R" "barry loves R"
```

Flow control means where the R brain goes in your code

```
good_mood <- FALSE

if(good_mood) {
  cat("Great!")
} else {
  cat("Bummer!")
}
```

Bummer!

A Protein generating function

Q. Write a function, called `generate_protein()` that generates a user specified length protein sequence.

There are 20 natural amino-acids

```
aa <- c("A", "R", "N", "D", "C", "Q", "E", "G", "H", "I", "L", "K", "M", "F", "P", "S", "T", "W")
```

```
generate_protein <- function(len) {
```

```

# The amino-acids to sample from
aa <- c("A", "R", "N", "D", "C", "Q", "E", "G", "H", "I", "L", "K", "M", "F", "P", "S", "T", "V", "W", "Y")
# Draw n=len amino-acids to make our sequence
ans <- sample(aa, size = len, replace=TRUE )
ans <- paste(ans, collapse ="")
return(ans)
}

```

```
my_seq <- generate_protein(42)  
my_seq
```

```
[1] "CTKCDWQWVHFVTVINSKCQAPTFDPKFPQNHGQKPSATSIN"
```

Q. Use that function to generate random protein sequences between length 6 and 12

```
generate_protein(6)
```

```
[1] "LDIAGC"
```

```
generate_protein(7)
```

```
[1] "AWFVSAS"
```

```
generate_protein(8)
```

```
[1] "DNDLVNTW"
```

```
generate_protein(9)
```

```
[1] "IFPPPRPDQ"
```

```
generate_protein(10)
```

```
[1] "ISTTSFITTE"
```

```
generate_protein(11)
```

```
[1] "TATCTNYERMV"
```

```
generate_protein(12)
```

```
[1] "HICTPMRVNQCY"
```

```
for(i in 6:12) {  
  #FASTA ID line ">id"  
  cat(">", i, sep = "", "\n")  
  # Protein sequence line  
  cat(generate_protein(i), "\n")  
}
```

```
>6  
WRELDL  
>7  
IAQHIYH  
>8
```

RFFICFPL
>9
FLTEAHLPE
>10
PSVYNMLQWR
>11
SGCPMGIFIQN
>12
DALGFTWIGNIK

Q. Are any of your sequences unique i.e. not found anywhere in nature?

My sequences that were 10, 11, and 12 amino-acids long were unique. These sequences did not have a 100% percent identity, and 100% coverage.