

## Chapter 3

# Data and Programming

### 3.1 Data Types

R has a number of basic data *types*.

- Numeric
  - Also known as Double. The default type when dealing with numbers.
  - Examples: 1, 1.0, 42.5
- Integer
  - Examples: 1L, 2L, 42L
- Complex
  - Example: 4 + 2i
- Logical
  - Two possible values: TRUE and FALSE
  - You can also use T and F, but this is *not* recommended.
  - NA is also considered logical.
- Character
  - Examples: "a", "Statistics", "1 plus 2."

### 3.2 Data Structures

R also has a number of basic data *structures*. A data structure is either homogeneous (all elements are of the same data type) or heterogeneous (elements can be of more than one data type).

Dimension	Homogeneous	Heterogeneous
1	Vector	List
2	Matrix	Data Frame
3+	Array	

### 3.2.1 Vectors

Many operations in R make heavy use of **vectors**. Vectors in R are indexed starting at 1. That is what the [1] in the output is indicating, that the first element of the row being displayed is the first element of the vector. Larger vectors will start additional rows with [\*] where \* is the index of the first element of the row.

Possibly the most common way to create a vector in R is using the `c()` function, which is short for “combine.” As the name suggests, it combines a list of elements separated by commas.

```
c(1, 3, 5, 7, 8, 9)
```

```
## [1] 1 3 5 7 8 9
```

Here R simply outputs this vector. If we would like to store this vector in a **variable** we can do so with the **assignment** operator `=`. In this case the variable `x` now holds the vector we just created, and we can access the vector by typing `x`.

```
x = c(1, 3, 5, 7, 8, 9)
x
```

```
## [1] 1 3 5 7 8 9
```

As an aside, there is a long history of the assignment operator in R, partially due to the keys available on the [keyboards of the creators of the S language](#). (Which preceded R.) For simplicity we will use `=`, but know that often you will see `<=` as the assignment operator.

The pros and cons of these two are well beyond the scope of this book, but know that for our purposes you will have no issue if you simply use `=`. If you are interested in the weird cases where the difference matters, check out [The R Inferno](#).

If you wish to use `<=`, you will still need to use `=`, however only for argument passing. Some users like to keep assignment (`<=`) and argument passing (`=`) separate. No matter what you choose, the more important thing is that you

**stay consistent.** Also, if working on a larger collaborative project, you should use whatever style is already in place.

Because vectors must contain elements that are all the same type, R will automatically coerce to a single type when attempting to create a vector that combines multiple types.

```
c(42, "Statistics", TRUE)
```

```
## [1] "42"          "Statistics" "TRUE"
```

```
c(42, TRUE)
```

```
## [1] 42 1
```

Frequently you may wish to create a vector based on a sequence of numbers. The quickest and easiest way to do this is with the `:` operator, which creates a sequence of integers between two specified integers.

```
(y = 1:100)
```

```
## [1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18
## [19] 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36
## [37] 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54
## [55] 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72
## [73] 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90
## [91] 91 92 93 94 95 96 97 98 99 100
```

Here we see R labeling the rows after the first since this is a large vector. Also, we see that by putting parentheses around the assignment, R both stores the vector in a variable called `y` and automatically outputs `y` to the console.

Note that scalars do not exist in R. They are simply vectors of length 1.

```
2
```

```
## [1] 2
```

If we want to create a sequence that isn't limited to integers and increasing by 1 at a time, we can use the `seq()` function.

```
seq(from = 1.5, to = 4.2, by = 0.1)
```

```
## [1] 1.5 1.6 1.7 1.8 1.9 2.0 2.1 2.2 2.3 2.4 2.5 2.6 2.7 2.8 2.9 3.0 3.1 3.2 3.3
## [20] 3.4 3.5 3.6 3.7 3.8 3.9 4.0 4.1 4.2
```

We will discuss functions in detail later, but note here that the input labels `from`, `to`, and `by` are optional.

```
seq(1.5, 4.2, 0.1)
```

```
## [1] 1.5 1.6 1.7 1.8 1.9 2.0 2.1 2.2 2.3 2.4 2.5 2.6 2.7 2.8 2.9 3.0 3.1 3.2 3.3
## [20] 3.4 3.5 3.6 3.7 3.8 3.9 4.0 4.1 4.2
```

Another common operation to create a vector is `rep()`, which can repeat a single value a number of times.

```
rep("A", times = 10)
```

```
## [1] "A" "A" "A" "A" "A" "A" "A" "A" "A" "A"
```

The `rep()` function can be used to repeat a vector some number of times.

```
rep(x, times = 3)
```

```
## [1] 1 3 5 7 8 9 1 3 5 7 8 9 1 3 5 7 8 9
```

We have now seen four different ways to create vectors:

- `c()`
- `:`
- `seq()`
- `rep()`

So far we have mostly used them in isolation, but they are often used together.

```
c(x, rep(seq(1, 9, 2), 3), c(1, 2, 3), 42, 2:4)
```

```
## [1] 1 3 5 7 8 9 1 3 5 7 9 1 3 5 7 9 1 3 5 7 9 1 2 3 42
## [26] 2 3 4
```

The length of a vector can be obtained with the `length()` function.

```
length(x)
```

```
## [1] 6
```

```
length(y)
```

```
## [1] 100
```

### 3.2.1.1 Subsetting

To subset a vector, we use square brackets, `[]`.

```
x
```

```
## [1] 1 3 5 7 8 9
```

```
x[1]
```

```
## [1] 1
```

```
x[3]
```

```
## [1] 5
```

We see that `x[1]` returns the first element, and `x[3]` returns the third element.

```
x[-2]
```

```
## [1] 1 5 7 8 9
```

We can also exclude certain indexes, in this case the second element.

```
x[1:3]
```

```
## [1] 1 3 5
```

```
x[c(1,3,4)]
```

```
## [1] 1 5 7
```

Lastly we see that we can subset based on a vector of indices.

All of the above are subsetting a vector using a vector of indexes. (Remember a single number is still a vector.) We could instead use a vector of logical values.

```
z = c(TRUE, TRUE, FALSE, TRUE, TRUE, FALSE)
z
```

```
## [1] TRUE TRUE FALSE TRUE TRUE FALSE
```

```
x[z]
```

```
## [1] 1 3 7 8
```

### 3.2.2 Vectorization

One of the biggest strengths of R is its use of vectorized operations. (Frequently the lack of understanding of this concept leads of a belief that R is *slow*. R is not the fastest language, but it has a reputation for being slower than it really is.)

```
x = 1:10
x + 1
```

```
## [1] 2 3 4 5 6 7 8 9 10 11
```

```
2 * x
```

```
## [1] 2 4 6 8 10 12 14 16 18 20
```

```
2 ^ x
```

```
## [1] 2 4 8 16 32 64 128 256 512 1024
```

```
sqrt(x)
```

```
## [1] 1.000000 1.414214 1.732051 2.000000 2.236068 2.449490 2.645751 2.828427
## [9] 3.000000 3.162278
```

```
log(x)
```

```
## [1] 0.0000000 0.6931472 1.0986123 1.3862944 1.6094379 1.7917595 1.9459101
## [8] 2.0794415 2.1972246 2.3025851
```

We see that when a function like `log()` is called on a vector `x`, a vector is returned which has applied the function to each element of the vector `x`.

### 3.2.3 Logical Operators

Operator	Summary	Example	Result
$x < y$	x less than y	$3 < 42$	TRUE
$x > y$	x greater than y	$3 > 42$	FALSE
$x \leq y$	x less than or equal to y	$3 \leq 42$	TRUE
$x \geq y$	x greater than or equal to y	$3 \geq 42$	FALSE
$x == y$	x equal to y	$3 == 42$	FALSE
$x != y$	x not equal to y	$3 != 42$	TRUE
$!x$	not x	$!(3 > 42)$	TRUE
$x   y$	x or y	$(3 > 42)   \text{TRUE}$	TRUE
$x \& y$	x and y	$(3 < 4) \& (42 > 13)$	TRUE

In R, logical operators are vectorized.

```
x = c(1, 3, 5, 7, 8, 9)
```

```
x > 3
```

```
## [1] FALSE FALSE TRUE TRUE TRUE TRUE
```

```
x < 3
```

```
## [1] TRUE FALSE FALSE FALSE FALSE FALSE
```

```
x == 3
```

```
## [1] FALSE TRUE FALSE FALSE FALSE FALSE
```

```
x != 3
```

```
## [1] TRUE FALSE TRUE TRUE TRUE TRUE
```

```
x == 3 & x != 3
```

```
## [1] FALSE FALSE FALSE FALSE FALSE FALSE
```

```
x == 3 | x != 3
```

```
## [1] TRUE TRUE TRUE TRUE TRUE TRUE
```

This is extremely useful for subsetting.

```
x[x > 3]
```

```
## [1] 5 7 8 9
```

```
x[x != 3]
```

```
## [1] 1 5 7 8 9
```

- TODO: coercion

```
sum(x > 3)
```

```
## [1] 4
```

```
as.numeric(x > 3)
```

```
## [1] 0 0 1 1 1 1
```

Here we see that using the `sum()` function on a vector of logical `TRUE` and `FALSE` values that is the result of `x > 3` results in a numeric result. R is first automatically coercing the logical to numeric where `TRUE` is 1 and `FALSE` is 0. This coercion from logical to numeric happens for most mathematical operations.

```
which(x > 3)
```

```
## [1] 3 4 5 6
```

```
x[which(x > 3)]
```

```
## [1] 5 7 8 9
```

```
max(x)
```

```
## [1] 9
```

```
which(x == max(x))
```

```
## [1] 6
```



```
which.max(x)
```

```
## [1] 6
```

### 3.2.4 More Vectorization

```
x = c(1, 3, 5, 7, 8, 9)
y = 1:100
```

```
x + 2
```

```
## [1] 3 5 7 9 10 11
```

```
x + rep(2, 6)
```

```
## [1] 3 5 7 9 10 11
```

```
x > 3
```

```
## [1] FALSE FALSE TRUE TRUE TRUE TRUE
```

```
x > rep(3, 6)
```

```
## [1] FALSE FALSE TRUE TRUE TRUE TRUE
```

```
x + y
```

```
## Warning in x + y: longer object length is not a multiple of shorter object
## length
```

```
## [1] 2 5 8 11 13 15 8 11 14 17 19 21 14 17 20 23 25 27
## [19] 20 23 26 29 31 33 26 29 32 35 37 39 32 35 38 41 43 45
## [37] 38 41 44 47 49 51 44 47 50 53 55 57 50 53 56 59 61 63
## [55] 56 59 62 65 67 69 62 65 68 71 73 75 68 71 74 77 79 81
## [73] 74 77 80 83 85 87 80 83 86 89 91 93 86 89 92 95 97 99
## [91] 92 95 98 101 103 105 98 101 104 107
```

```
length(x)
```

```
## [1] 6
```

```
length(y)
```

```
## [1] 100
```

```
length(y) / length(x)
```

```
## [1] 16.66667
```

```
(x + y) - y
```

```
## Warning in x + y: longer object length is not a multiple of shorter object
## length
```

```
## [1] 1 3 5 7 8 9 1 3 5 7 8 9 1 3 5 7 8 9 1 3 5 7 8 9 1 3 5 7 8 9 1
## [38] 3 5 7 8 9 1 3 5 7 8 9 1 3 5 7 8 9 1 3 5 7 8 9 1 3 5 7 8 9 1 3
## [75] 5 7 8 9 1 3 5 7 8 9 1 3 5 7 8 9 1 3 5 7 8 9 1 3 5 7
```

```
y = 1:60
x + y
```

```
## [1] 2 5 8 11 13 15 8 11 14 17 19 21 14 17 20 23 25 27 20 23 26 29 31 33 26
## [26] 29 32 35 37 39 32 35 38 41 43 45 38 41 44 47 49 51 44 47 50 53 55 57 50 53
## [51] 56 59 61 63 56 59 62 65 67 69
```

```
length(y) / length(x)
```

```
## [1] 10
```

```
rep(x, 10) + y
```

```
## [1] 2 5 8 11 13 15 8 11 14 17 19 21 14 17 20 23 25 27 20 23 26 29 31 33 26
## [26] 29 32 35 37 39 32 35 38 41 43 45 38 41 44 47 49 51 44 47 50 53 55 57 50 53
## [51] 56 59 61 63 56 59 62 65 67 69
```

```
all(x + y == rep(x, 10) + y)
```

```
## [1] TRUE
```

```
identical(x + y, rep(x, 10) + y)
```

```
## [1] TRUE
```

```
# ?any
# ?all.equal
```

### 3.2.5 Matrices

R can also be used for **matrix** calculations. Matrices have rows and columns containing a single data type. In a matrix, the order of rows and columns is important. (This is not true of *data frames*, which we will see later.)

Matrices can be created using the **matrix** function.

```
x = 1:9
x
```

```
## [1] 1 2 3 4 5 6 7 8 9
```

```
X = matrix(x, nrow = 3, ncol = 3)
X
```

```
##      [,1] [,2] [,3]
## [1,]    1    4    7
## [2,]    2    5    8
## [3,]    3    6    9
```

Note here that we are using two different variables: lower case **x**, which stores a vector and capital **X**, which stores a matrix. (Following the usual mathematical convention.) We can do this because R is case sensitive.

By default the **matrix** function reorders a vector into columns, but we can also tell R to use rows instead.

```
Y = matrix(x, nrow = 3, ncol = 3, byrow = TRUE)
Y
```

```
##      [,1] [,2] [,3]
## [1,]    1    2    3
## [2,]    4    5    6
## [3,]    7    8    9
```

We can also create a matrix of a specified dimension where every element is the same, in this case 0.

```
Z = matrix(0, 2, 4)
Z
```

```
##      [,1] [,2] [,3] [,4]
## [1,]    0    0    0    0
## [2,]    0    0    0    0
```

Like vectors, matrices can be subsetted using square brackets, `[]`. However, since matrices are two-dimensional, we need to specify both a row and a column when subsetting.

```
X
```

```
##      [,1] [,2] [,3]
## [1,]    1    4    7
## [2,]    2    5    8
## [3,]    3    6    9
```

```
X[1, 2]
```

```
## [1] 4
```

Here we accessed the element in the first row and the second column. We could also subset an entire row or column.

```
X[1, ]
```

```
## [1] 1 4 7
```

```
X[, 2]
```

```
## [1] 4 5 6
```

We can also use vectors to subset more than one row or column at a time. Here we subset to the first and third column of the second row.

```
X[2, c(1, 3)]
```

```
## [1] 2 8
```

Matrices can also be created by combining vectors as columns, using `cbind`, or combining vectors as rows, using `rbind`.

```
x = 1:9
rev(x)
```

```
## [1] 9 8 7 6 5 4 3 2 1
```

```
rep(1, 9)
```

```
## [1] 1 1 1 1 1 1 1 1 1
```

```
rbind(x, rev(x), rep(1, 9))
```

```
##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9]
## x      1    2    3    4    5    6    7    8    9
##      9    8    7    6    5    4    3    2    1
##      1    1    1    1    1    1    1    1    1
```

```
cbind(col_1 = x, col_2 = rev(x), col_3 = rep(1, 9))
```

```
##      col_1 col_2 col_3
## [1,]     1     9     1
## [2,]     2     8     1
## [3,]     3     7     1
## [4,]     4     6     1
## [5,]     5     5     1
## [6,]     6     4     1
## [7,]     7     3     1
## [8,]     8     2     1
## [9,]     9     1     1
```

When using `rbind` and `cbind` you can specify “argument” names that will be used as column names.

R can then be used to perform matrix calculations.

```
x = 1:9
y = 9:1
X = matrix(x, 3, 3)
Y = matrix(y, 3, 3)
X
```

```
##      [,1] [,2] [,3]
## [1,]    1    4    7
## [2,]    2    5    8
## [3,]    3    6    9
```

```
Y
```

```
##      [,1] [,2] [,3]
## [1,]    9    6    3
## [2,]    8    5    2
## [3,]    7    4    1
```

```
X + Y
```

```
##      [,1] [,2] [,3]
## [1,]   10   10   10
## [2,]   10   10   10
## [3,]   10   10   10
```

```
X - Y
```

```
##      [,1] [,2] [,3]
## [1,]   -8   -2    4
## [2,]   -6    0    6
## [3,]   -4    2    8
```

```
X * Y
```

```
##      [,1] [,2] [,3]
## [1,]    9   24   21
## [2,]   16   25   16
## [3,]   21   24    9
```

```
X / Y
```

```
##           [,1]      [,2]      [,3]
## [1,] 0.1111111 0.6666667 2.333333
## [2,] 0.2500000 1.0000000 4.000000
## [3,] 0.4285714 1.5000000 9.000000
```

Note that `X * Y` is not matrix multiplication. It is element by element multiplication. (Same for `X / Y`). Instead, matrix multiplication uses `%*%`. Other matrix functions include `t()` which gives the transpose of a matrix and `solve()` which returns the inverse of a square matrix if it is invertible.

```
X %*% Y
```

```
##           [,1] [,2] [,3]
## [1,]      90   54   18
## [2,]     114   69   24
## [3,]     138   84   30
```

```
t(X)
```

```
##           [,1] [,2] [,3]
## [1,]        1    2    3
## [2,]        4    5    6
## [3,]        7    8    9
```

```
Z = matrix(c(9, 2, -3, 2, 4, -2, -3, -2, 16), 3, byrow = TRUE)
Z
```

```
##           [,1] [,2] [,3]
## [1,]        9    2   -3
## [2,]        2    4   -2
## [3,]       -3   -2   16
```

```
solve(Z)
```

```
##           [,1]      [,2]      [,3]
## [1,] 0.12931034 -0.05603448 0.01724138
## [2,] -0.05603448 0.29094828 0.02586207
## [3,] 0.01724138 0.02586207 0.06896552
```

To verify that `solve(Z)` returns the inverse, we multiply it by `Z`. We would expect this to return the identity matrix, however we see that this is not the case due to some computational issues. However, R also has the `all.equal()` function which checks for equality, with some small tolerance which accounts for some computational issues. The `identical()` function is used to check for exact equality.

```
solve(Z) %*% Z
```

```
##           [,1]      [,2]      [,3]
## [1,] 1.000000e+00 -6.245005e-17 0.000000e+00
## [2,] 8.326673e-17  1.000000e+00 5.551115e-17
## [3,] 2.775558e-17  0.000000e+00 1.000000e+00
```

```
diag(3)
```

```
##           [,1] [,2] [,3]
## [1,]      1    0    0
## [2,]      0    1    0
## [3,]      0    0    1
```

```
all.equal(solve(Z) %*% Z, diag(3))
```

```
## [1] TRUE
```

R has a number of matrix specific functions for obtaining dimension and summary information.

```
X = matrix(1:6, 2, 3)
X
```

```
##           [,1] [,2] [,3]
## [1,]      1    3    5
## [2,]      2    4    6
```

```
dim(X)
```

```
## [1] 2 3
```

```
rowSums(X)
```

```
## [1]  9 12
```

```
colSums(X)
```

```
## [1]  3  7 11
```



```
rowMeans(X)
```

```
## [1] 3 4
```

```
colMeans(X)
```

```
## [1] 1.5 3.5 5.5
```

The `diag()` function can be used in a number of ways. We can extract the diagonal of a matrix.

```
diag(Z)
```

```
## [1] 9 4 16
```

Or create a matrix with specified elements on the diagonal. (And 0 on the off-diagonals.)

```
diag(1:5)
```

```
##      [,1] [,2] [,3] [,4] [,5]
## [1,] 1    0    0    0    0
## [2,] 0    2    0    0    0
## [3,] 0    0    3    0    0
## [4,] 0    0    0    4    0
## [5,] 0    0    0    0    5
```

Or, lastly, create a square matrix of a certain dimension with 1 for every element of the diagonal and 0 for the off-diagonals.

```
diag(5)
```

```
##      [,1] [,2] [,3] [,4] [,5]
## [1,] 1    0    0    0    0
## [2,] 0    1    0    0    0
## [3,] 0    0    1    0    0
## [4,] 0    0    0    1    0
## [5,] 0    0    0    0    1
```

### Calculations with Vectors and Matrices

Certain operations in R, for example `%%` have different behavior on vectors and matrices. To illustrate this, we will first create two vectors.

```
a_vec = c(1, 2, 3)
b_vec = c(2, 2, 2)
```

Note that these are indeed vectors. They are not matrices.

```
c(is.vector(a_vec), is.vector(b_vec))
```

```
## [1] TRUE TRUE
```

```
c(is.matrix(a_vec), is.matrix(b_vec))
```

```
## [1] FALSE FALSE
```

When this is the case, the `%%` operator is used to calculate the **dot product**, also know as the **inner product** of the two vectors.

The dot product of vectors  $a = [a_1, a_2, \dots, a_n]$  and  $b = [b_1, b_2, \dots, b_n]$  is defined to be

$$a \cdot b = \sum_{i=1}^n a_i b_i = a_1 b_1 + a_2 b_2 + \dots a_n b_n.$$

```
a_vec %% b_vec # inner product
```

```
##      [,1]
## [1,]    12
```

```
a_vec %o% b_vec # outer product
```

```
##      [,1] [,2] [,3]
## [1,]    2    2    2
## [2,]    4    4    4
## [3,]    6    6    6
```

The `%o%` operator is used to calculate the **outer product** of the two vectors.

When vectors are coerced to become matrices, they are column vectors. So a vector of length  $n$  becomes an  $n \times 1$  matrix after coercion.

```
as.matrix(a_vec)
```

```
##      [,1]
## [1,]    1
## [2,]    2
## [3,]    3
```

If we use the `%%` operator on matrices, `%%` again performs the expected matrix multiplication. So you might expect the following to produce an error, because the dimensions are incorrect.

```
as.matrix(a_vec) %% b_vec
```

```
##      [,1] [,2] [,3]
## [1,]    2    2    2
## [2,]    4    4    4
## [3,]    6    6    6
```

At face value this is a  $3 \times 1$  matrix, multiplied by a  $3 \times 1$  matrix. However, when `b_vec` is automatically coerced to be a matrix, R decided to make it a “row vector”, a  $1 \times 3$  matrix, so that the multiplication has conformable dimensions.

If we had coerced both, then R would produce an error.

```
as.matrix(a_vec) %% as.matrix(b_vec)
```

Another way to calculate a *dot product* is with the `crossprod()` function. Given two vectors, the `crossprod()` function calculates their dot product. The function has a rather misleading name.

```
crossprod(a_vec, b_vec) # inner product
```

```
##      [,1]
## [1,]   12
```

```
tcrossprod(a_vec, b_vec) # outer product
```

```
##      [,1] [,2] [,3]
## [1,]    2    2    2
## [2,]    4    4    4
## [3,]    6    6    6
```

These functions could be very useful later. When used with matrices  $X$  and  $Y$  as arguments, it calculates

$$X^T Y.$$

When dealing with linear models, the calculation

$$X^T X$$

is used repeatedly.

```
C_mat = matrix(c(1, 2, 3, 4, 5, 6), 2, 3)
D_mat = matrix(c(2, 2, 2, 2, 2, 2), 2, 3)
```

This is useful both as a shortcut for a frequent calculation and as a more efficient implementation than using `t()` and `%*%`.

```
crossprod(C_mat, D_mat)
```

```
##      [,1] [,2] [,3]
## [1,]    6    6    6
## [2,]   14   14   14
## [3,]   22   22   22
```

```
t(C_mat) %*% D_mat
```

```
##      [,1] [,2] [,3]
## [1,]    6    6    6
## [2,]   14   14   14
## [3,]   22   22   22
```

```
all.equal(crossprod(C_mat, D_mat), t(C_mat) %*% D_mat)
```

```
## [1] TRUE
```

```
crossprod(C_mat, C_mat)
```

```
##      [,1] [,2] [,3]
## [1,]    5   11   17
## [2,]   11   25   39
## [3,]   17   39   61
```

```
t(C_mat) %*% C_mat
```

```
##      [,1] [,2] [,3]
## [1,]    5   11   17
## [2,]   11   25   39
## [3,]   17   39   61
```

```
all.equal(crossprod(C_mat, C_mat), t(C_mat) %*% C_mat)
```

```
## [1] TRUE
```

### 3.2.6 Lists

A list is a one-dimensional heterogeneous data structure. So it is indexed like a vector with a single integer value, but each element can contain an element of any type.

```
# creation
list(42, "Hello", TRUE)
```

```
## [[1]]
## [1] 42
##
## [[2]]
## [1] "Hello"
##
## [[3]]
## [1] TRUE
```

```
ex_list = list(
  a = c(1, 2, 3, 4),
  b = TRUE,
  c = "Hello!",
  d = function(arg = 42) {print("Hello World!")},
  e = diag(5)
)
```

Lists can be subset using two syntaxes, the `$` operator, and square brackets `[]`. The `$` operator returns a named **element** of a list. The `[]` syntax returns a **list**, while the `[[ ]]` returns an **element** of a list.

- `ex_list[1]` returns a list containing the first element.

- `ex_list[[1]]` returns the first element of the list, in this case, a vector.

```
# subsetting
ex_list$e
```

```
##      [,1] [,2] [,3] [,4] [,5]
## [1,]    1    0    0    0    0
## [2,]    0    1    0    0    0
## [3,]    0    0    1    0    0
## [4,]    0    0    0    1    0
## [5,]    0    0    0    0    1
```

```
ex_list[1:2]
```

```
## $a
## [1] 1 2 3 4
##
## $b
## [1] TRUE
```

```
ex_list[1]
```

```
## $a
## [1] 1 2 3 4
```

```
ex_list[[1]]
```

```
## [1] 1 2 3 4
```

```
ex_list[c("e", "a")]
```

```
## $e
##      [,1] [,2] [,3] [,4] [,5]
## [1,]    1    0    0    0    0
## [2,]    0    1    0    0    0
## [3,]    0    0    1    0    0
## [4,]    0    0    0    1    0
## [5,]    0    0    0    0    1
##
## $a
## [1] 1 2 3 4
```

```
ex_list["e"]
```

```
## $e
##      [,1] [,2] [,3] [,4] [,5]
## [1,]    1    0    0    0    0
## [2,]    0    1    0    0    0
## [3,]    0    0    1    0    0
## [4,]    0    0    0    1    0
## [5,]    0    0    0    0    1
```

```
ex_list[["e"]]
```

```
##      [,1] [,2] [,3] [,4] [,5]
## [1,]    1    0    0    0    0
## [2,]    0    1    0    0    0
## [3,]    0    0    1    0    0
## [4,]    0    0    0    1    0
## [5,]    0    0    0    0    1
```

```
ex_list$d
```

```
## function(arg = 42) {print("Hello World!")}
```

```
ex_list$d(arg = 1)
```

```
## [1] "Hello World!"
```

### 3.2.7 Data Frames

We have previously seen vectors and matrices for storing data as we introduced R. We will now introduce a **data frame** which will be the most common way that we store and interact with data in this course.

```
example_data = data.frame(x = c(1, 3, 5, 7, 9, 1, 3, 5, 7, 9),
                          y = c(rep("Hello", 9), "Goodbye"),
                          z = rep(c(TRUE, FALSE), 5))
```

Unlike a matrix, which can be thought of as a vector rearranged into rows and columns, a data frame is not required to have the same data type for each element. A data frame is a **list** of vectors. So, each vector must contain the same data type, but the different vectors can store different data types.

```
example_data
```

```
##      x      y      z
## 1  1  Hello TRUE
## 2  3  Hello FALSE
## 3  5  Hello TRUE
## 4  7  Hello FALSE
## 5  9  Hello TRUE
## 6  1  Hello FALSE
## 7  3  Hello TRUE
## 8  5  Hello FALSE
## 9  7  Hello TRUE
## 10 9 Goodbye FALSE
```

Unlike a list which has more flexibility, the elements of a data frame must all be vectors, and have the same length.

```
example_data$x
```

```
## [1] 1 3 5 7 9 1 3 5 7 9
```

```
all.equal(length(example_data$x),
          length(example_data$y),
          length(example_data$z))
```

```
## [1] TRUE
```

```
str(example_data)
```

```
## 'data.frame': 10 obs. of 3 variables:
## $ x: num 1 3 5 7 9 1 3 5 7 9
## $ y: chr "Hello" "Hello" "Hello" "Hello" ...
## $ z: logi TRUE FALSE TRUE FALSE TRUE FALSE ...
```

```
nrow(example_data)
```

```
## [1] 10
```

```
ncol(example_data)
```

```
## [1] 3
```



```
dim(example_data)
```

```
## [1] 10 3
```

The `data.frame()` function above is one way to create a data frame. We can also import data from various file types into R, as well as use data stored in packages.

The example data above can also be found here as a [.csv file](#). To read this data into R, we would use the `read_csv()` function from the `readr` package. Note that R has a built in function `read.csv()` that operates very similarly. The `readr` function `read_csv()` has a number of advantages. For example, it is much faster reading larger data. [It also uses the `tibble` package to read the data as a tibble.](#)

```
library(readr)
example_data_from_csv = read_csv("data/example-data.csv")
```

This particular line of code assumes that the file `example_data.csv` exists in a folder called `data` in your current working directory.

```
example_data_from_csv
```

```
## # A tibble: 10 x 3
##       x y      z
##   <dbl> <chr> <lgl>
## 1     1 1 Hello  TRUE
## 2     2 3 Hello  FALSE
## 3     3 5 Hello  TRUE
## 4     4 7 Hello  FALSE
## 5     5 9 Hello  TRUE
## 6     6 1 Hello  FALSE
## 7     7 3 Hello  TRUE
## 8     8 5 Hello  FALSE
## 9     9 7 Hello  TRUE
## 10    9 Goodbye FALSE
```

A tibble is simply a data frame that prints with sanity. Notice in the output above that we are given additional information such as dimension and variable type.

The `as_tibble()` function can be used to coerce a regular data frame to a tibble.

```
library(tibble)
example_data = as_tibble(example_data)
example_data
```

```
## # A tibble: 10 x 3
##       x y      z
##   <dbl> <chr> <lgl>
## 1     1 1 Hello  TRUE
## 2     3 3 Hello  FALSE
## 3     5 5 Hello  TRUE
## 4     7 7 Hello  FALSE
## 5     9 9 Hello  TRUE
## 6     1 1 Hello  FALSE
## 7     3 3 Hello  TRUE
## 8     5 5 Hello  FALSE
## 9     7 7 Hello  TRUE
## 10    9 9 Goodbye FALSE
```

Alternatively, we could use the “Import Dataset” feature in RStudio which can be found in the environment window. (By default, the top-right pane of RStudio.) Once completed, this process will automatically generate the code to import a file. The resulting code will be shown in the console window. In recent versions of RStudio, `read_csv()` is used by default, thus reading in a tibble.

Earlier we looked at installing packages, in particular the `ggplot2` package. (A package for visualization. While not necessary for this course, it is quickly growing in popularity.)

```
library(ggplot2)
```

Inside the `ggplot2` package is a dataset called `mpg`. By loading the package using the `library()` function, we can now access `mpg`.

When using data from inside a package, there are three things we would generally like to do:

- Look at the raw data.
- Understand the data. (Where did it come from? What are the variables? Etc.)
- Visualize the data.

To look at the data, we have two useful commands: `head()` and `str()`.

```
head(mpg, n = 10)
```

```
## # A tibble: 10 x 11
##   manufacturer model    displ  year   cyl trans  drv    cty   hwy fl    class
##   <chr>          <chr>    <dbl> <int> <int> <chr>  <chr> <int> <int> <chr> <chr>
## 1 audi          a4         1.8  1999     4 auto(l~ f      18    29 p    comp~
## 2 audi          a4         1.8  1999     4 manual~ f      21    29 p    comp~
## 3 audi          a4         2    2008     4 manual~ f      20    31 p    comp~
## 4 audi          a4         2    2008     4 auto(a~ f      21    30 p    comp~
## 5 audi          a4         2.8  1999     6 auto(l~ f      16    26 p    comp~
## 6 audi          a4         2.8  1999     6 manual~ f      18    26 p    comp~
## 7 audi          a4         3.1  2008     6 auto(a~ f      18    27 p    comp~
## 8 audi          a4 quat~ 1.8  1999     4 manual~ 4      18    26 p    comp~
## 9 audi          a4 quat~ 1.8  1999     4 auto(l~ 4      16    25 p    comp~
## 10 audi         a4 quat~ 2    2008     4 manual~ 4      20    28 p    comp~
```

The function `head()` will display the first `n` observations of the data frame. The `head()` function was more useful before tibbles. Notice that `mpg` is a tibble already, so the output from `head()` indicates there are only 10 observations. Note that this applies to `head(mpg, n = 10)` and not `mpg` itself. Also note that tibbles print a limited number of rows and columns by default. The last line of the printed output indicates which rows and columns were omitted.

```
mpg
```

```
## # A tibble: 234 x 11
##   manufacturer model    displ  year   cyl trans  drv    cty   hwy fl    class
##   <chr>          <chr>    <dbl> <int> <int> <chr>  <chr> <int> <int> <chr> <chr>
## 1 audi          a4         1.8  1999     4 auto(l~ f      18    29 p    comp~
## 2 audi          a4         1.8  1999     4 manual~ f      21    29 p    comp~
## 3 audi          a4         2    2008     4 manual~ f      20    31 p    comp~
## 4 audi          a4         2    2008     4 auto(a~ f      21    30 p    comp~
## 5 audi          a4         2.8  1999     6 auto(l~ f      16    26 p    comp~
## 6 audi          a4         2.8  1999     6 manual~ f      18    26 p    comp~
## 7 audi          a4         3.1  2008     6 auto(a~ f      18    27 p    comp~
## 8 audi          a4 quat~ 1.8  1999     4 manual~ 4      18    26 p    comp~
## 9 audi          a4 quat~ 1.8  1999     4 auto(l~ 4      16    25 p    comp~
## 10 audi         a4 quat~ 2    2008     4 manual~ 4      20    28 p    comp~
## # ... with 224 more rows
```

The function `str()` will display the “structure” of the data frame. It will display the number of **observations** and **variables**, list the variables, give the type of each variable, and show some elements of each variable. This information can also be found in the “Environment” window in RStudio.

```
str(mpg)
```

```
## tibble [234 x 11] (S3: tbl_df/tbl/data.frame)
## $ manufacturer: chr [1:234] "audi" "audi" "audi" "audi" ...
## $ model       : chr [1:234] "a4" "a4" "a4" "a4" ...
## $ displ       : num [1:234] 1.8 1.8 2 2 2.8 2.8 3.1 1.8 1.8 2 ...
## $ year        : int [1:234] 1999 1999 2008 2008 1999 1999 2008 1999 1999 2008 ...
## $ cyl         : int [1:234] 4 4 4 4 6 6 6 4 4 4 ...
## $ trans       : chr [1:234] "auto(l5)" "manual(m5)" "manual(m6)" "auto(av)" ...
## $ drv         : chr [1:234] "f" "f" "f" "f" ...
## $ cty         : int [1:234] 18 21 20 21 16 18 18 18 16 20 ...
## $ hwy         : int [1:234] 29 29 31 30 26 26 27 26 25 28 ...
## $ fl         : chr [1:234] "p" "p" "p" "p" ...
## $ class       : chr [1:234] "compact" "compact" "compact" "compact" ...
```

It is important to note that while matrices have rows and columns, data frames (tibbles) instead have observations and variables. When displayed in the console or viewer, each row is an observation and each column is a variable. However generally speaking, their order does not matter, it is simply a side-effect of how the data was entered or stored.

In this dataset an observation is for a particular model-year of a car, and the variables describe attributes of the car, for example its highway fuel efficiency.

To understand more about the data set, we use the `?` operator to pull up the documentation for the data.

```
?mpg
```

R has a number of functions for quickly working with and extracting basic information from data frames. To quickly obtain a vector of the variable names, we use the `names()` function.

```
names(mpg)
```

```
## [1] "manufacturer" "model"          "displ"          "year"          "cyl"
## [6] "trans"        "drv"            "cty"           "hwy"           "fl"
## [11] "class"
```

To access one of the variables **as a vector**, we use the `$` operator.

```
mpg$year
```

```
## [1] 1999 1999 2008 2008 1999 1999 2008 1999 1999 2008 2008 1999 1999 2008 2008
## [16] 1999 2008 2008 2008 2008 1999 2008 1999 2008 1999 1999 2008 2008 2008 2008
## [31] 1999 1999 1999 2008 1999 2008 2008 1999 1999 1999 1999 2008 2008 2008 1999
## [46] 1999 2008 2008 2008 2008 1999 2008 1999 2008 2008 2008 1999 1999 1999 2008
## [61] 2008 1999 2008 1999 2008 1999 2008 2008 2008 2008 1999 1999 2008 1999 1999
## [76] 1999 2008 1999 1999 1999 2008 2008 1999 1999 1999 1999 2008 1999 2008
## [91] 1999 1999 2008 2008 1999 1999 2008 2008 2008 1999 1999 1999 1999 1999 2008
## [106] 2008 2008 2008 1999 1999 2008 2008 1999 1999 2008 1999 1999 2008 2008 2008
## [121] 2008 2008 2008 2008 1999 1999 2008 2008 2008 2008 1999 2008 2008 1999 1999
## [136] 1999 2008 1999 2008 2008 1999 1999 1999 2008 2008 2008 2008 1999 1999 2008
## [151] 1999 1999 2008 2008 1999 1999 1999 2008 2008 1999 1999 2008 2008 2008 2008
## [166] 1999 1999 1999 1999 2008 2008 2008 2008 1999 1999 1999 2008 2008 1999
## [181] 1999 2008 2008 1999 1999 2008 1999 1999 2008 1999 1999 1999 2008 1999
## [196] 1999 2008 2008 1999 2008 1999 2008 1999 2008 2008 2008 1999 1999 2008
## [211] 2008 1999 1999 1999 1999 2008 2008 2008 1999 1999 1999 1999 1999 1999
## [226] 2008 2008 1999 1999 2008 2008 1999 1999 2008
```

```
mpg$hwy
```

```
## [1] 29 29 31 30 26 26 27 26 25 28 27 25 25 25 25 24 25 23 20 15 20 17 17 26 23
## [26] 26 25 24 19 14 15 17 27 30 26 29 26 24 24 22 22 24 24 17 22 21 23 23 19 18
## [51] 17 17 19 19 12 17 15 17 17 12 17 16 18 15 16 12 17 17 16 12 15 16 17 15 17
## [76] 17 18 17 19 17 19 19 17 17 17 16 16 17 15 17 26 25 26 24 21 22 23 22 20 33
## [101] 32 32 29 32 34 36 36 29 26 27 30 31 26 26 28 26 29 28 27 24 24 24 22 19 20
## [126] 17 12 19 18 14 15 18 18 15 17 16 18 17 19 19 17 29 27 31 32 27 26 26 25 25
## [151] 17 17 20 18 26 26 27 28 25 25 24 27 25 26 23 26 26 26 26 25 27 25 27 20 20
## [176] 19 17 20 17 29 27 31 31 26 26 28 27 29 31 31 26 26 27 30 33 35 37 35 15 18
## [201] 20 20 22 17 19 18 20 29 26 29 29 24 44 29 26 29 29 29 29 23 24 44 41 29 26
## [226] 28 29 29 29 28 29 26 26 26
```

We can use the `dim()`, `nrow()` and `ncol()` functions to obtain information about the dimension of the data frame.

```
dim(mpg)
```

```
## [1] 234 11
```

```
nrow(mpg)
```

```
## [1] 234
```

```
ncol(mpg)
```

```
## [1] 11
```

Here `nrow()` is also the number of observations, which in most cases is the *sample size*.

Subsetting data frames can work much like subsetting matrices using square brackets, `[,]`. Here, we find fuel efficient vehicles earning over 35 miles per gallon and only display `manufacturer`, `model` and `year`.

```
mpg[mpg$hwy > 35, c("manufacturer", "model", "year")]
```

```
## # A tibble: 6 x 3
##   manufacturer model      year
##   <chr>         <chr>    <int>
## 1 honda         civic     2008
## 2 honda         civic     2008
## 3 toyota        corolla   2008
## 4 volkswagen    jetta     1999
## 5 volkswagen    new beetle 1999
## 6 volkswagen    new beetle 1999
```

An alternative would be to use the `subset()` function, which has a much more readable syntax.

```
subset(mpg, subset = hwy > 35, select = c("manufacturer", "model", "year"))
```

Lastly, we could use the `filter` and `select` functions from the `dplyr` package which introduces the `%>%` operator from the `magrittr` package. This is not necessary for this course, however the `dplyr` package is something you should be aware of as it is becoming a popular tool in the R world.

```
library(dplyr)
mpg %>% filter(hwy > 35) %>% select(manufacturer, model, year)
```

All three approaches produce the same results. Which you use will be largely based on a given situation as well as user preference.

When subsetting a data frame, be aware of what is being returned, as sometimes it may be a vector instead of a data frame. Also note that there are differences between subsetting a data frame and a tibble. A data frame operates more like a matrix where it is possible to reduce the subset to a vector. A tibble operates more like a list where it always subsets to another tibble.

## 3.3 Programming Basics

### 3.3.1 Control Flow

In R, the if/else syntax is:

```
if (...) {  
  some R code  
} else {  
  more R code  
}
```

For example,

```
x = 1  
y = 3  
if (x > y) {  
  z = x * y  
  print("x is larger than y")  
} else {  
  z = x + 5 * y  
  print("x is less than or equal to y")  
}
```

```
## [1] "x is less than or equal to y"
```

```
z
```

```
## [1] 16
```

R also has a special function `ifelse()` which is very useful. It returns one of two specified values based on a conditional statement.

```
ifelse(4 > 3, 1, 0)
```

```
## [1] 1
```

The real power of `ifelse()` comes from its ability to be applied to vectors.

```
fib = c(1, 1, 2, 3, 5, 8, 13, 21)  
ifelse(fib > 6, "Foo", "Bar")
```

```
## [1] "Bar" "Bar" "Bar" "Bar" "Bar" "Foo" "Foo" "Foo"
```

Now a `for` loop example,

```
x = 11:15
for (i in 1:5) {
  x[i] = x[i] * 2
}

x
```

```
## [1] 22 24 26 28 30
```

Note that this `for` loop is very normal in many programming languages, but not in R. In R we would not use a loop, instead we would simply use a vectorized operation.

```
x = 11:15
x = x * 2
x
```

```
## [1] 22 24 26 28 30
```

### 3.3.2 Functions

So far we have been using functions, but haven't actually discussed some of their details.

```
function_name(arg1 = 10, arg2 = 20)
```

To use a function, you simply type its name, followed by an open parenthesis, then specify values of its arguments, then finish with a closing parenthesis.

An **argument** is a variable which is used in the body of the function. Specifying the values of the arguments is essentially providing the inputs to the function.

We can also write our own functions in R. For example, we often like to “standardize” variables, that is, subtracting the sample mean, and dividing by the sample standard deviation.

$$\frac{x - \bar{x}}{s}$$

In R we would write a function to do this. When writing a function, there are three things you must do.



- Give the function a name. Preferably something that is short, but descriptive.
- Specify the arguments using `function()`
- Write the body of the function within curly braces, `{}`.

```
standardize = function(x) {
  m = mean(x)
  std = sd(x)
  result = (x - m) / std
  result
}
```

Here the name of the function is `standardize`, and the function has a single argument `x` which is used in the body of function. Note that the output of the final line of the body is what is returned by the function. In this case the function returns the vector stored in the variable `result`.

To test our function, we will take a random sample of size `n = 10` from a normal distribution with a mean of 2 and a standard deviation of 5.

```
(test_sample = rnorm(n = 10, mean = 2, sd = 5))
```

```
## [1] -0.1456927 -5.5692281  2.9219036  8.6167755  2.8873106 -3.1477341
## [7]  0.1278323 -0.5568065  1.5051560  1.9603510
```

```
standardize(x = test_sample)
```

```
## [1] -0.2634915 -1.6844768  0.5402294  2.0323058  0.5311659 -1.0500369
## [7] -0.1918270 -0.3712048  0.1690366  0.2882993
```

This function could be written much more succinctly, simply performing all the operations on one line and immediately returning the result, without storing any of the intermediate results.

```
standardize = function(x) {
  (x - mean(x)) / sd(x)
}
```

When specifying arguments, you can provide default arguments.

```
power_of_num = function(num, power = 2) {
  num ^ power
}
```

Let's look at a number of ways that we could run this function to perform the operation  $10^2$  resulting in 100.

```
power_of_num(10)
```

```
## [1] 100
```

```
power_of_num(10, 2)
```

```
## [1] 100
```

```
power_of_num(num = 10, power = 2)
```

```
## [1] 100
```

```
power_of_num(power = 2, num = 10)
```

```
## [1] 100
```

Note that without using the argument names, the order matters. The following code will not evaluate to the same output as the previous example.

```
power_of_num(2, 10)
```

```
## [1] 1024
```

Also, the following line of code would produce an error since arguments without a default value must be specified.

```
power_of_num(power = 5)
```

To further illustrate a function with a default argument, we will write a function that calculates sample variance two ways.

By default, it will calculate the unbiased estimate of  $\sigma^2$ , which we will call  $s^2$ .

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x - \bar{x})^2$$

It will also have the ability to return the biased estimate (based on maximum likelihood) which we will call  $\hat{\sigma}^2$ .

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x - \bar{x})^2$$

```
get_var = function(x, biased = FALSE) {  
  n = length(x) - 1 * !biased  
  (1 / n) * sum((x - mean(x)) ^ 2)  
}
```

```
get_var(test_sample)
```

```
## [1] 14.56753
```

```
get_var(test_sample, biased = FALSE)
```

```
## [1] 14.56753
```

```
var(test_sample)
```

```
## [1] 14.56753
```

We see the function is working as expected, and when returning the unbiased estimate it matches R's built in function `var()`. Finally, let's examine the biased estimate of  $\sigma^2$ .

```
get_var(test_sample, biased = TRUE)
```

```
## [1] 13.11077
```

