# 02-05: Vector Operations and NA

## 1 - Purpose

- mathematical operations on vectors
- dealing with missing values in a vector
- R shortcuts for vector mathematics

## 2 - Questions about the material...

If you have any questions about the material in this lesson or your Class Project *feel free to email them to the instructor here*.

## 3 - Vector math

For this lesson, we are going to use the data from the file **LansingWeather2.csv**.  First we need to open the CSV file and save the data in the file to a data frame, which we will call **weatherData**.

```
1  {
2    rm(list=ls());  options(show.error.locations = TRUE);
3
4    # read data from LansingWeather2.csv and save to the variable weatherData...
5    weatherData = read.csv("data/LansingWeather2.csv");
6  }
```

The data frame, **weatherData**, has the date, high temperature, low temperature, and precipitation for 14 days (March 27-April 9, 2017) in Lansing, MI.
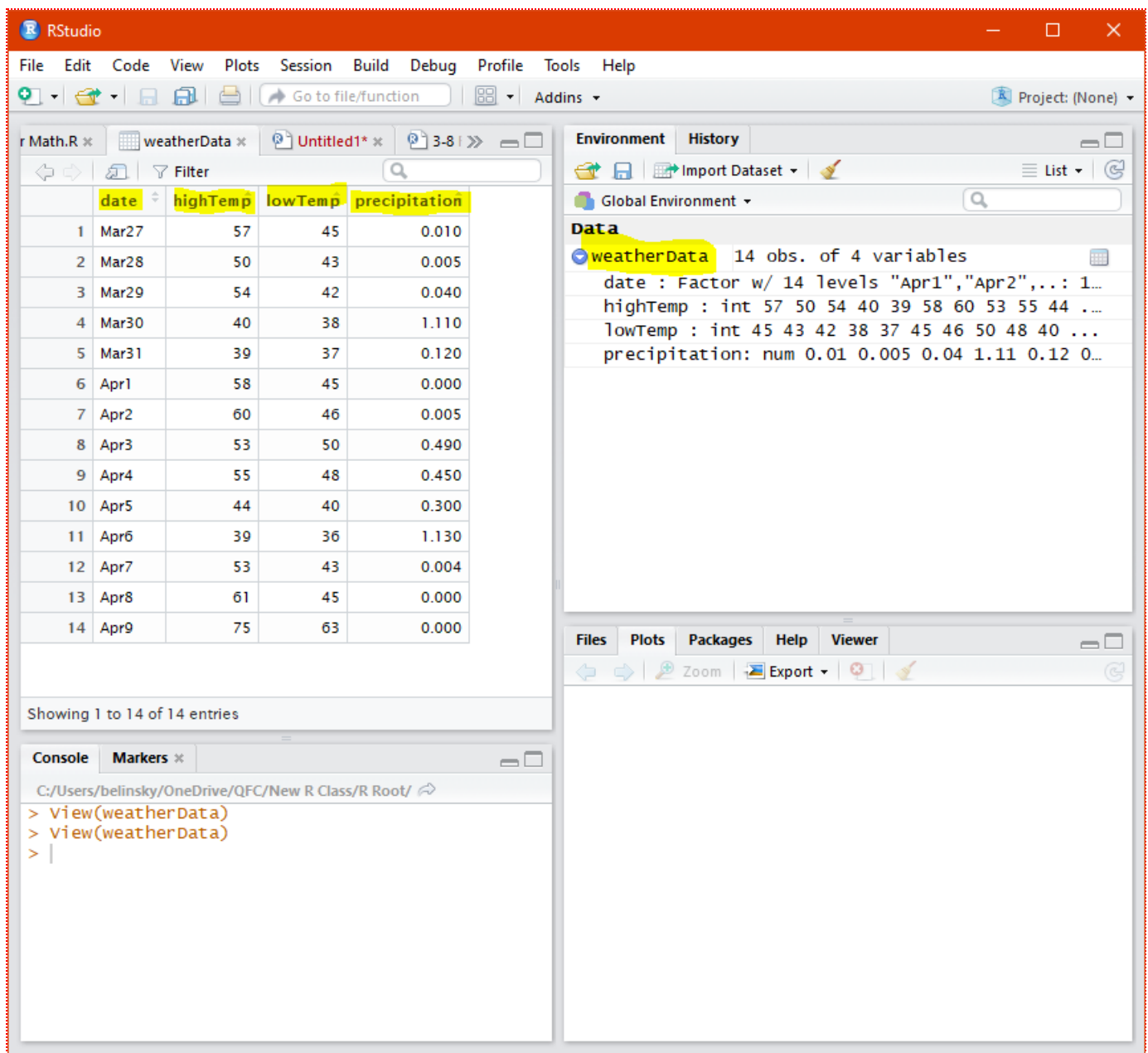
*Fig 1: Looking at the values in the **weatherData** data frame in the Main Window.*

## 3.1 - Finding change in temperature

We want to find the temperature range for each day -- in other words, the high temperature minus the low temperature. This means subtracting each day's low temperature from the high temperature and saving that value to a vector.

First we will extract the low and high temperature values from ***weatherData*** and save them to a vector. There are two ways we can do this:
1) by column number

```
1  # get all values from the 3rd column in weatherData
2  highTemp = weatherData[, 3];
```

2) or by column name

```
1  # get all values from the column named "lowTemp" in weatherData
2  lowTemp = weatherData[, "lowTemp"];
```

Note: the row number is left blank to indicate we are getting values *from all rows.*

We will use a *for()* to iterate through each value (i.e., the highTemp - lowTemp for each day) and solve for the temperature difference.  To do this we also need a vector that will hold the change in temperature values, which we will call *changeInTemp*.

```
1  changeInTemp = c();   #declared a vector
```

*changeInTemp* acts as a state variable because:
1) *changeInTemp* gets initialized before the *for()*
2) *changeInTemp* gets populated during the iterations of the *for()*
3) *changeInTemp* final state is the full set of temperature differences

## 3.2 - Subtracting values from two different vectors

We are going to use a *for()* loop that iterates through each value in *highTemp* and *lowTemp*, finds the difference, and saves the answer to *changeInTemp*.  We need to know the length of the vectors to do this.  In this case, we already know there are **14** values, but we will use *length()* to calculate this value so that it works for vectors of any size (i.e., if we wanted to run the same code but use more days).

```
1  vectorLength = length(lowTemp);
```

Note: since the length of all columns in a data frame are by definition the same, we only need to get the length of one vector.

We use *vectorLength* to create a sequence that the *for()* iterates through.  The *for()* iterates through the sequence **1:vectorLength**, **(1:14)**, and assigns the value of the **14 highTemp** - **lowTemp** difference operations to the **14 changeInTemp** vector elements.

```
1  for(i in 1:vectorLength)
2  {
3     changeInTemp[i] = highTemp[i] - lowTemp[i];
4  }
```

Let's put all this code together:

```
1  {
2     rm(list=ls()); options(show.error.locations = TRUE);
3
4     # read data from LansingWeather2.csv and save to variable weatherData
5     weatherData = read.csv("data/LansingWeather2.csv");
6     # get all values from the 3rd column in weatherData
```

```
 7    highTemp = weatherData[, 2];
 8    # get all values from the column named "lowTemp" in weatherData
 9    lowTemp = weatherData[, "lowTemp"];
10
11    changeInTemp = c(); # declare a vector
12    vectorLength = length(lowTemp); # vectorLength will be 14 (length of data)
13
14    # go through the sequence 1:14
15    for(i in 1:vectorLength)
16    {
17        # subtract lowTemp from highTemp for all 14 values and save to changeInTemp
18        changeInTemp[i] = highTemp[i] - lowTemp[i];
19    }
20 }
```

The Environment Window shows the vector **changeInTemp** but it only shows up to 10 values.  If you want to see all the values in **changeInTemp**, you can type *changeInTemp* in the Console Window and the 14 values in **changeInTemp** will appear on the next line.
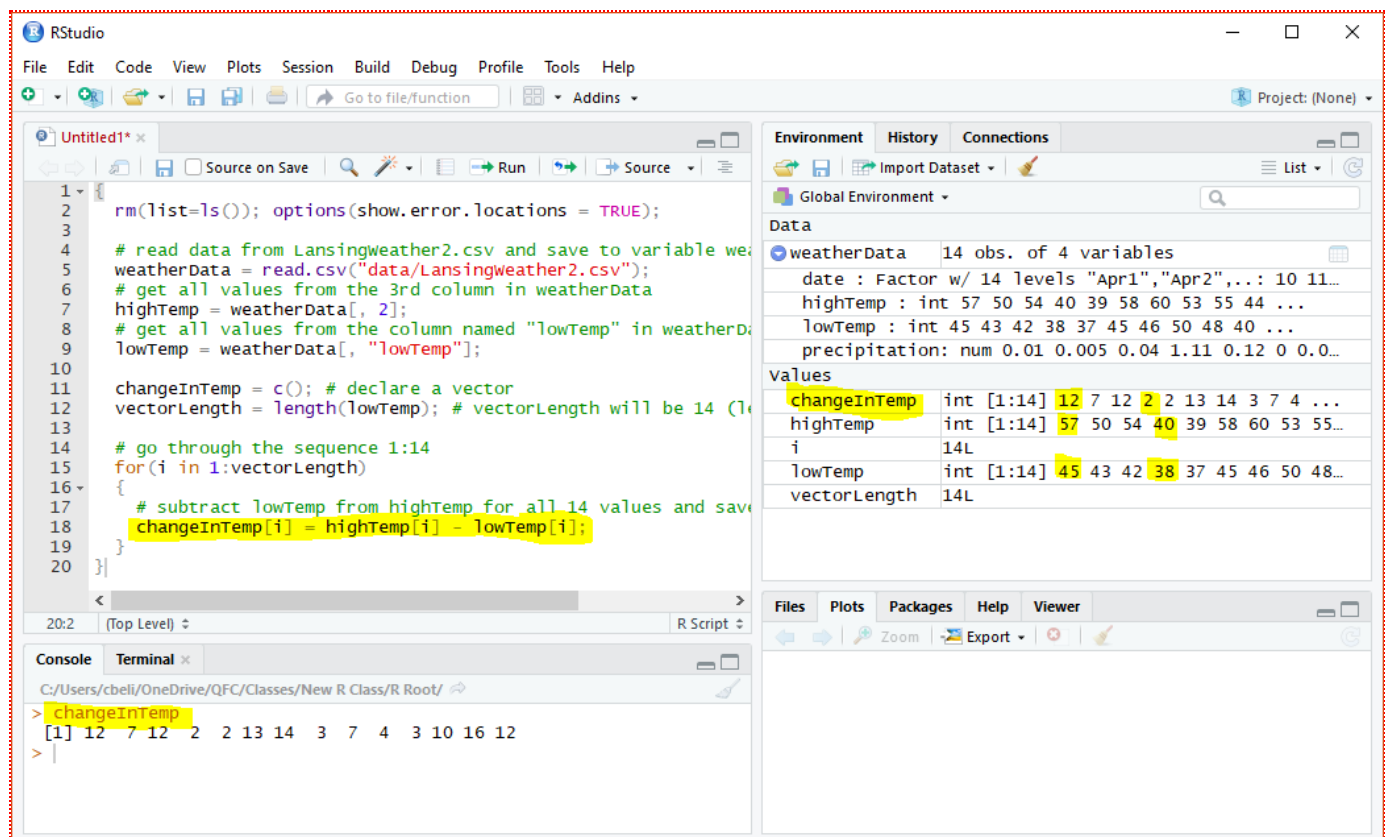


*Fig 2: Using a **for()** to perform iterative mathematical operations on vectors.*

# 4 - Dealing with missing values

In the previous example, we have an idealized situation where every day had a high and low temperature associated with it.  In the real world, especially with large amounts of data, there are often missing data.  In R,

missing data are represented by *NA,* but the CSV file could designate missing data according to a number of different conventions.

First lets create a data set with missing values and save this as *MissingTemps.csv* in the *Data* directory

```
 1  date,highTemp, lowTemp, precipitation
 2  Mar27,57,45,0.01
 3  Mar28,50,43,0.005
 4  Mar29,54, ,0.04
 5  Mar30,40,38,1.11
 6  Mar31,39,NA,0.12
 7  Apr1,58,45,0
 8  Apr2,60, ,0.005
 9  Apr3,53,50,0.49
10  Apr4,55,48,0.45
11  Apr5,44,40,0.30
12  Apr6,39,36,1.13
13  Apr7,NULL,43,0.004
14  Apr8,61,45,0
15  Apr9,75,63,0
```

In the above data set, the *lowTemp* for **Mar29** and **Apr2** are left blank, the *lowTemp* for **Mar31** is given as *NA,* and the *highTemp* for **Apr7** is given as *NULL*.  So, there are four missing values in *missingTemps.csv*.

## 4.1 - Viewing blank or NA values

We are going to open a new script and save *missingTemps.csv* data to a Data Frame called *weatherData,* but we want to standardize the way that missing values are recorded.  So, in *read.csv()* we add the parameter *na.strings*:

```
 1  na.strings=c("", " ", "NULL", NULL, "NA", "na", NA, "null")
```

Essentially, *na.strings* is a vector that contains all the value that you want R to assign as *NA*.  It is good to be paranoid here and think of all the possible ways in which CSV files will present these values!  In this case, we only needed to use " ", "NULL", and "NA".

```
 1  {
 2      rm(list=ls()); options(show.error.locations = TRUE);
 3
 4      weatherData = read.csv("data/missingTemps.csv",
 5              na.strings=c("", " ", "NULL", NULL, "NA", "na", NA, "null"));
 6  }
```

Execute the code above and double-click on **weatherData** in the **Environment Window** so that it appears in the Main Window. We see that all four missing values in **weatherData** are now labelled as **NA**.
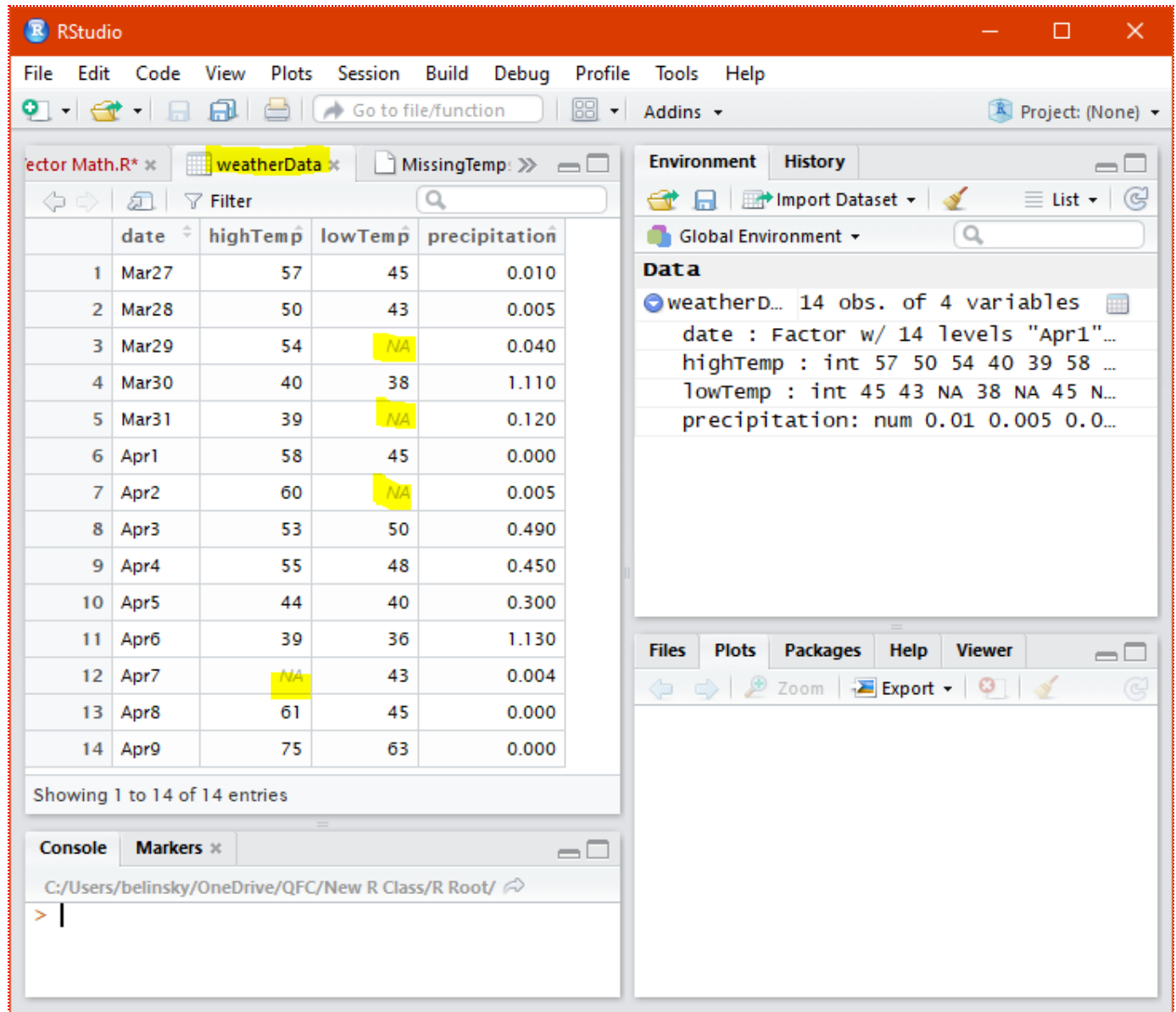


Fig 3: **NA** values in a data frame

## 4.2 - Mathematical operations on vectors with NA values

We are going to run the same change in temperature script as before except with NA values in the data frame:

```
1  {
2    rm(list=ls()); options(show.error.locations = TRUE);
3
4    # read data from missingTemps.csv and save to the variable weatherData...
5    weatherData = read.csv("data/missingTemps.csv",
6        na.strings=c("", " ", "NULL", NULL, "NA", "na", NA, "null"));
7
```

```
 8    # get all values from the 3rd column in weatherData
 9    highTemp = weatherData[, 2];
10    # get all values from the column named "lowTemp" in weatherData
11    lowTemp = weatherData[, "lowTemp"];
12
13    changeInTemp = c(); # declare a vector
14    vectorLength = length(lowTemp); # vectorLength will be 14 (length of data)
15
16    # go through the sequence 1:14
17    for(i in 1:vectorLength)
18    {
19        # subtract lowTemp from highTemp for all 14 values and save to changeInTemp
20        changeInTemp[i] = highTemp[i] - lowTemp[i];
21    }
22 }
```

Note: anytime there is an **NA** in a calculation (e.g., 3rd element in **changeInTemp**) , the answer is going to be **NA**.  This is a special feature of **NA --** if you used any other value to represent missing data then you would get an error in the calculation.



*Fig 4: Doing mathematical operation with **NA** values*

*Extension: What counts as NA*

# 5 - Vector Operation Shortcuts

So far we have used *for()* loops to perform iterative operations on a vector or multiple vectors. However, R has built in functions that do all these operations using a lot less code. For example, we can perform the above change in temperature calculation without a *for()*.

```r
1  {
2     rm(list=ls()); options(show.error.locations = TRUE);
3
4     weatherData = read.csv("data/missingTemps.csv",
5             na.strings=c("", " ", "NULL", NULL, "NA", "na", NA, "null"));
6
7     highTemp = weatherData[, 2];
8     lowTemp = weatherData[, "lowTemp"];
9
10    changeInTemp = highTemp - lowTemp;
11 }
```

Line 9 does all the work of iterating through the values in the vector, subtracting the values, and saving the answer to the vector *changeInTemp*. The result is the same as the script that uses the *for()* loop.
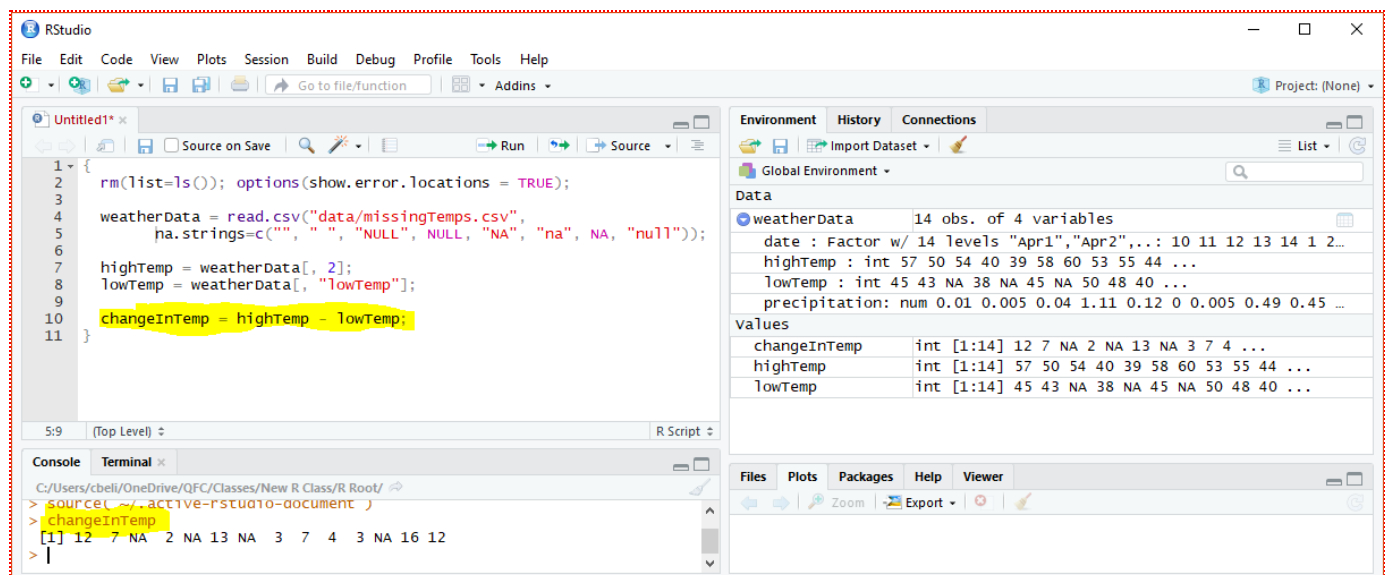


*Fig 5: Subtracting a vector in R.*

This method is obviously easier than *for()* loops for subtracting two vectors -- and it is the method R programmer will usually use. However, this method is really a shortcut and it is important to understand what is going on in with the *for()* example because, eventually, you will find a situation where you need to use a *for()* to solve a more complex problem.

## 5.1 - Many other shortcuts

R has many functions that can quickly perform the most common operations on vectors like *sum(), max(), min(),* and *mean()*.

One important parameter used in all of these functions is *na.rm:*

***na.rm = TRUE*** tells R to exclude the ***NA*** values from the vector before performing the operation.
***na.rm = FALSE*** tells R to return ***NA*** if there are ***NA*** values in the vector.

The default value for ***na.rm*** is ***FALSE***.

```
1  {
2    rm(list=ls()); options(show.error.locations = TRUE);
3
4    weatherData = read.csv("data/missingTemps.csv",
5            na.strings=c("", " ", "NULL", NULL, "NA", "na", NA, "null"))
6
7    # save high temp and low temp columns to vectors
8    highTemp = weatherData[, 2];
9    lowTemp = weatherData[, "lowTemp"];
10
11   # get min and max temp (na.rm = default = FALSE)
12   minTemp = min(lowTemp);
13   maxTemp = max(highTemp);
14
15   # get min and max temp ignoring the NAs
16   minTempTake2 = min(lowTemp, na.rm=TRUE);
17   maxTempTake2 = max(highTemp, na.rm=TRUE);
18
19   # get the index of min and max values
20   minIndex = which.min(lowTemp);
21   maxIndex = which.max(highTemp);
22
23   # get simple statistics on a vector
24   sumTemp = sum(highTemp, na.rm=TRUE);
25   meanTemp = mean(highTemp, na.rm=TRUE);
26   medianTemp = median(highTemp, na.rm=TRUE);
27   stanDev = sd(highTemp, na.rm=TRUE);
28   variance = var(highTemp, na.rm=TRUE);
29
30   # get index of values that meet a condition
31   # this is more advanced and will be covered in detail later
32   whichHighGT60 = which(highTemp > 60);
33   whichLowLT43 = which(lowTemp < 43);
34  }
```
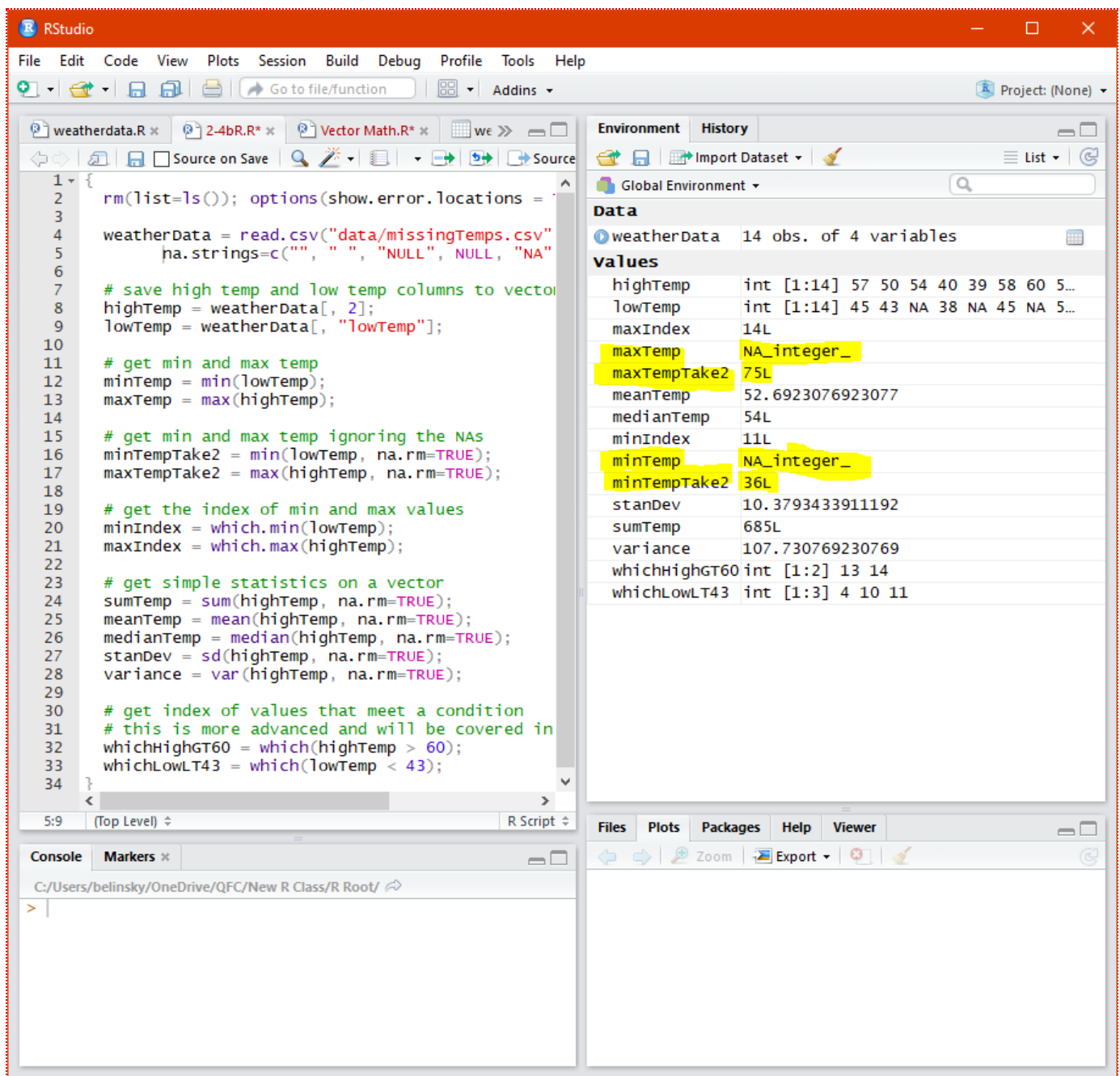
Fig 6: Using R functions to perform common vector operations.

# 6 - Application

*If you have any questions regarding this application or your Class Project, feel free to email them to the instructor here.  You can attach the whole Root Folder as a zipped file.*

1) Create a vector called ***changeInHighTemp*** and, using the data from **lansingWeather2.csv**, find the change in high temperatures from day-to-day.

     So, if the four high temperatures are: **40, 45, 35, 42**
     then the three ***changeInHighTemp*** would be: (45-40), (35-45), and (42-35) = **5, -10, 7**

2) There will be one less value in ***changeInHighTemp*** than in ***lansingWeather2.csv***.  We want the length of ***changeInHighTemp*** to be the same so we can add it to the data frame (remember, all vectors must be the

same size in a data frame).  To do this, add an **NA** value as the first value in ***changeInHighTemp***.  The **NA** indicates that there was no previous value to subtract and, hence, this value is not applicable.

3) Add the vector ***changeInHighTemp*** (with the **NA** as the first value) to the data frame.

*Save you script file as **app2-5.r**  in the **scripts** folder of your RStudio Project for the class.*

# 7 - Extension: What counts as NA

R is somewhat savvy when it moves data from a CSV to a data frame.  If there is an empty value in a column that is all numbers, then R will convert the empty value to **NA**.   **NULL**, however, is not considered to be a valid data frame value and if there is a **NULL** in a column, it will be treated as the string value "NULL".  R also does not consider lowercase **na** to be **NA** and this will also be treated as the string value "na".

There is now a second problem: if there is even one string in a column, all values in the column are treated as strings.  This means that **40** will be seen as the string "40" and mathematical operations on strings will cause an error.  Thus, specifying the ***na.strings*** argument can be useful to ensure that you will nor have to convert columns to numeric using an extra step after the data are imported.