

03-03 Matrices

1 - Purpose

- Shape weather data into a two-dimensional structure (i.e., has rows and columns)
- Perform mathematics and simple statistics on the rows, columns, and whole structure
- Write the matrix to a file

2 - Concepts

3 - Script for the lesson

[The script for this lesson is located here](#). We will be going through the script in order. You might want to first comment out all the lines and uncomment the lines as you go through the lesson.

If the NOAA/NCDC website is down, [you will need to download this rdata file](#), which contains the weather data for this lesson. Save this file to your **data** folder in the **R Root** directory. The script file for this lesson has instructions on how to adapt the script for the rdata file.

4 - Two-dimensional data

A data frame is a two dimensional data structure that is a collection of related vectors; for instance different types of weather data collected at the same times, as shown in the last lesson.

Data frames, by design, have *columns that hold different types of data*. This means that mathematical or statistical operations can only be performed on (numeric) columns as it would not make sense to perform these operation across row (e.g., taking the mean of mixed precipitation and temperature values or columns that have string values like precipitation type).

However, we can make a two-dimensional structure where all of the values are the same type, called a **matrix**. In a matrix, you can perform mathematical and statistical operations across rows, down columns, or on the entire structure. The matrix we will use in this lesson contains temperatures values for all days in January over six years (2011 to 2016) (*Fig 1*).

	Jan 2011	Jan 2012	Jan 2013	Jan 2014	Jan 2015	Jan 2016
1	54	44	25	14	29	29
2	25	29	26	13	34	35
3	33	23	29	13	33	33
...
30	26	29	57	19	32	29
31	24	41	56	28	23	47

Fig 1: Partial matrix of January temperatures from 2011 through 2016 that we will create in this lesson

5 - The matrix

In the previous lesson, we got weather data from NOAA/NCDC and put it into a data frame. A data frame is a 2D data structure that contains multiple vectors of equal length. In this case the vectors represent different weather measurements like high temperature, wind speed, or precipitation. The rows of a data frame represent the instances of each vector -- in the case of our weather data, the rows represent days.

In this lesson, we will introduce the *matrix*, which is like a data frame except that all columns have the same type of value. A data frame can have factors or strings like weather conditions (e.g., snow, fog, or rain) in one column and numerical values like temperatures in another column whereas *all values in a matrix must be of one type*.

One advantage to putting data in a matrix is that R can perform statistical operations faster on a matrix than on a data frame -- this becomes noticeable as your data set gets larger or as you iterate procedures multiple times. Another advantage is that R can perform statistical and mathematical operation on rows of a matrix, columns of a matrix, or the whole matrix.

6 - Getting temperature data

We are going to get data from the NOAA/NCDC database. However, this time we are only getting one piece of data, the max temperature (**TMAX**) from all days in January, 2011.

```
1 lansweather11 = ncdc(datasetid="GHCND",
2                       datatypeid=c("TMAX"),
3                       stationid="GHCND:USW00014836",
4                       startdate = "2011-01-01", enddate = "2011-01-31",
5                       token=myToken,
6                       limit=50 );
```

NOAA/NCDC returns a list object and we just want the data frame contained within the list. So, we need to subset the data frame from the list.

```
1 lansweather11Data = lansweather11[["data"]];
```

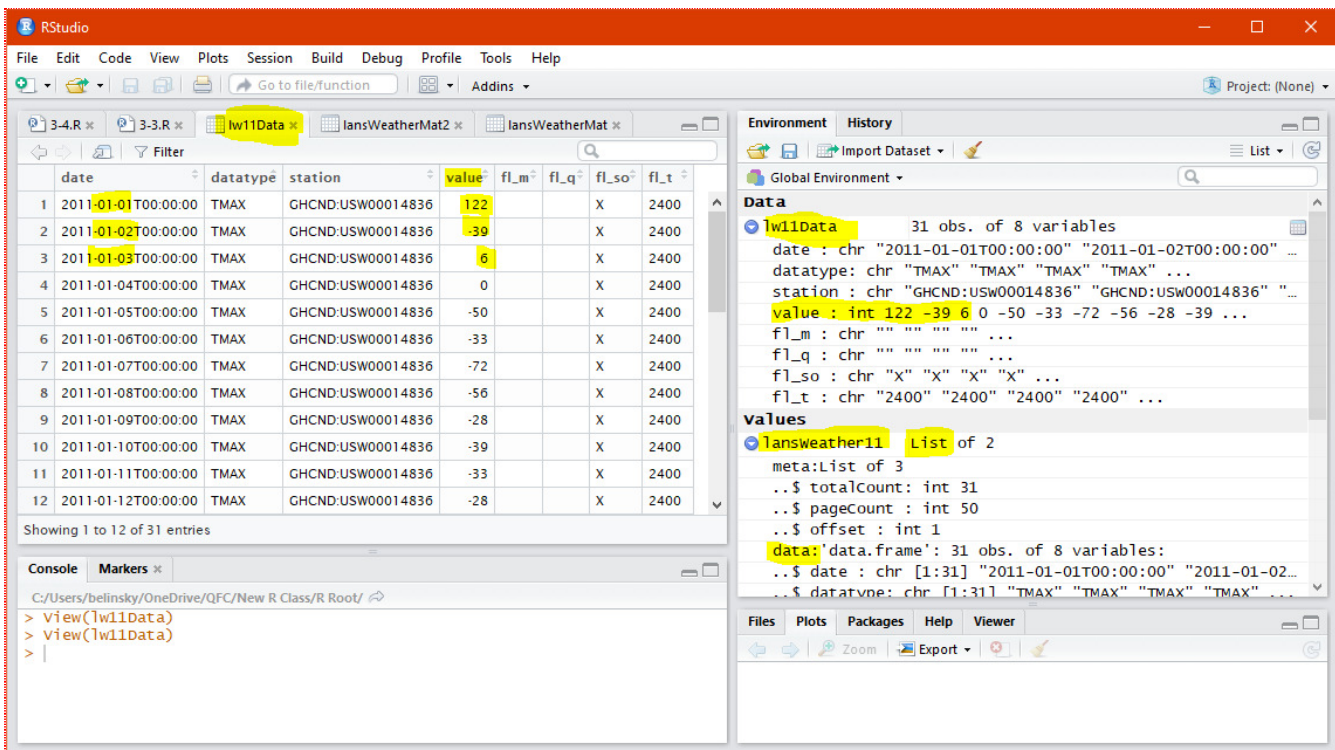


Fig 2: Getting the data from the list object returned by NOAA/NCDC's database

Because we only have one data type, the data frame is much simpler than last lesson. In this case, each row has a unique date and a maximum temperature (remember the temperature is in tenths of a Celsius degree -- we will take care of that in a bit).

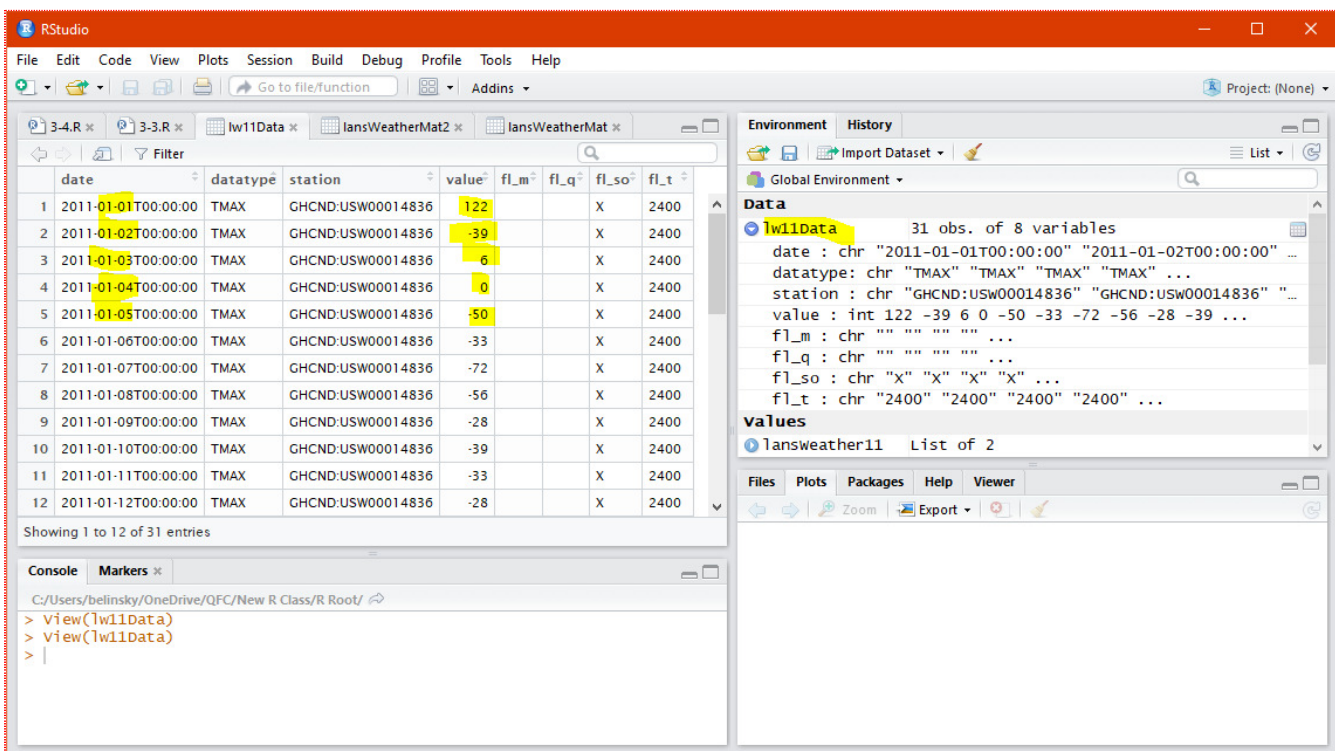


Fig 3: Data Frame showing unique date for each row and an associated temperature.

The **value** column holds the maximum temperatures, so we need to save the **value** column to a vector:

```
1 | lansweather11val = lansweather11Data[, "value"];
```

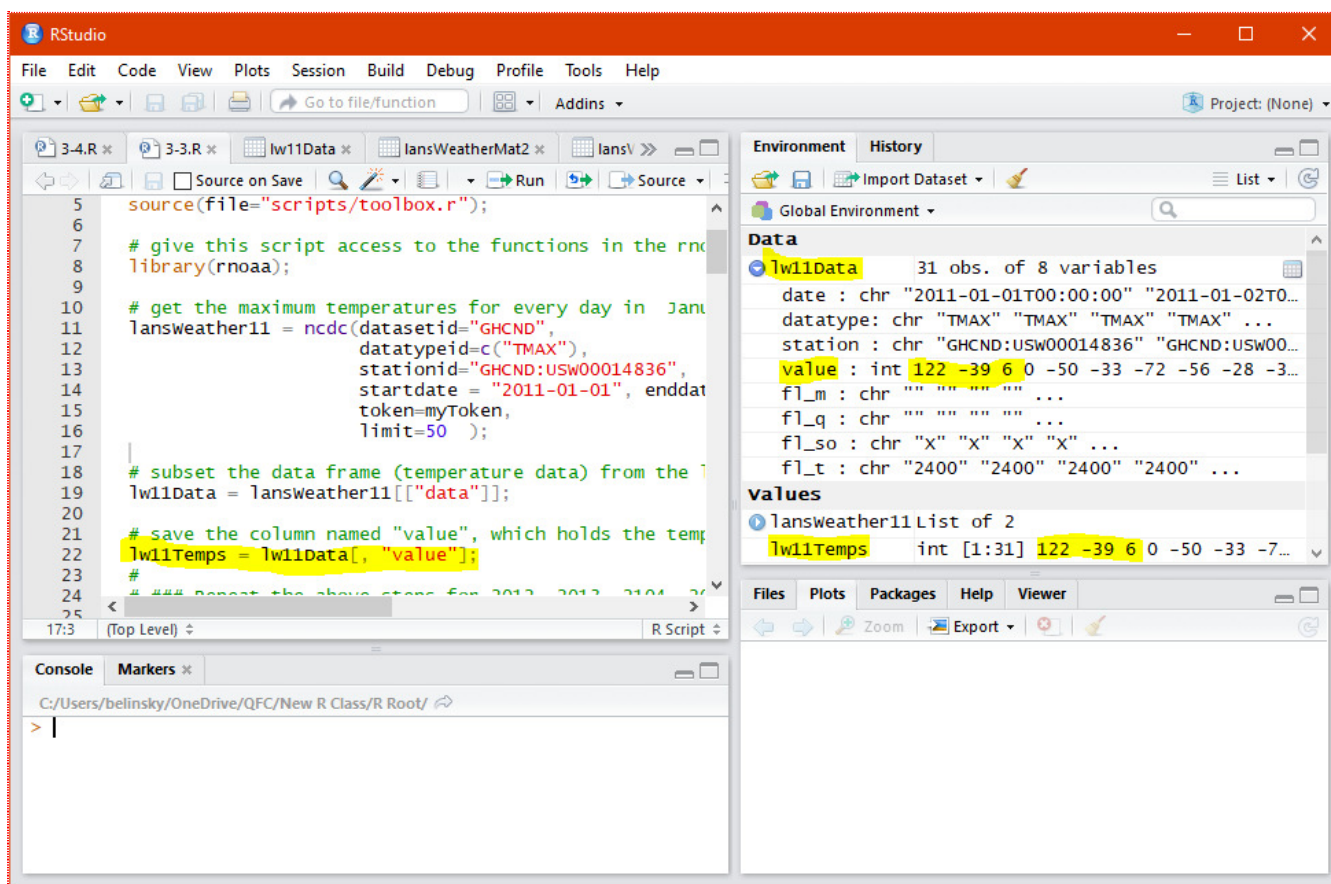


Fig 4: Saving the "value" column from the data frame to a vector

6.1 - Repeat for all six years

Now we are going to repeat the code above for the years 2012, 2013, 2014, 2015, and 2016. Note: **X** represents the last digit for the years 2012 to 2016.

```

1 lansweather1X = ncdc(datasetid="GHCND",
2                      datatypeid=c("TMX"),
3                      stationid="GHCND:USW00014836",
4                      startdate = "201X-01-01", enddate = "201X-01-31",
5                      token=myToken,
6                      limit=50 );

```

Sometimes the database does not accept multiple requests at once so we add a delay before each request.

We use **sys.sleep()** to add a delay, with the parameter **time = 0.5**.

```

1 sys.sleep(time=0.5);

```

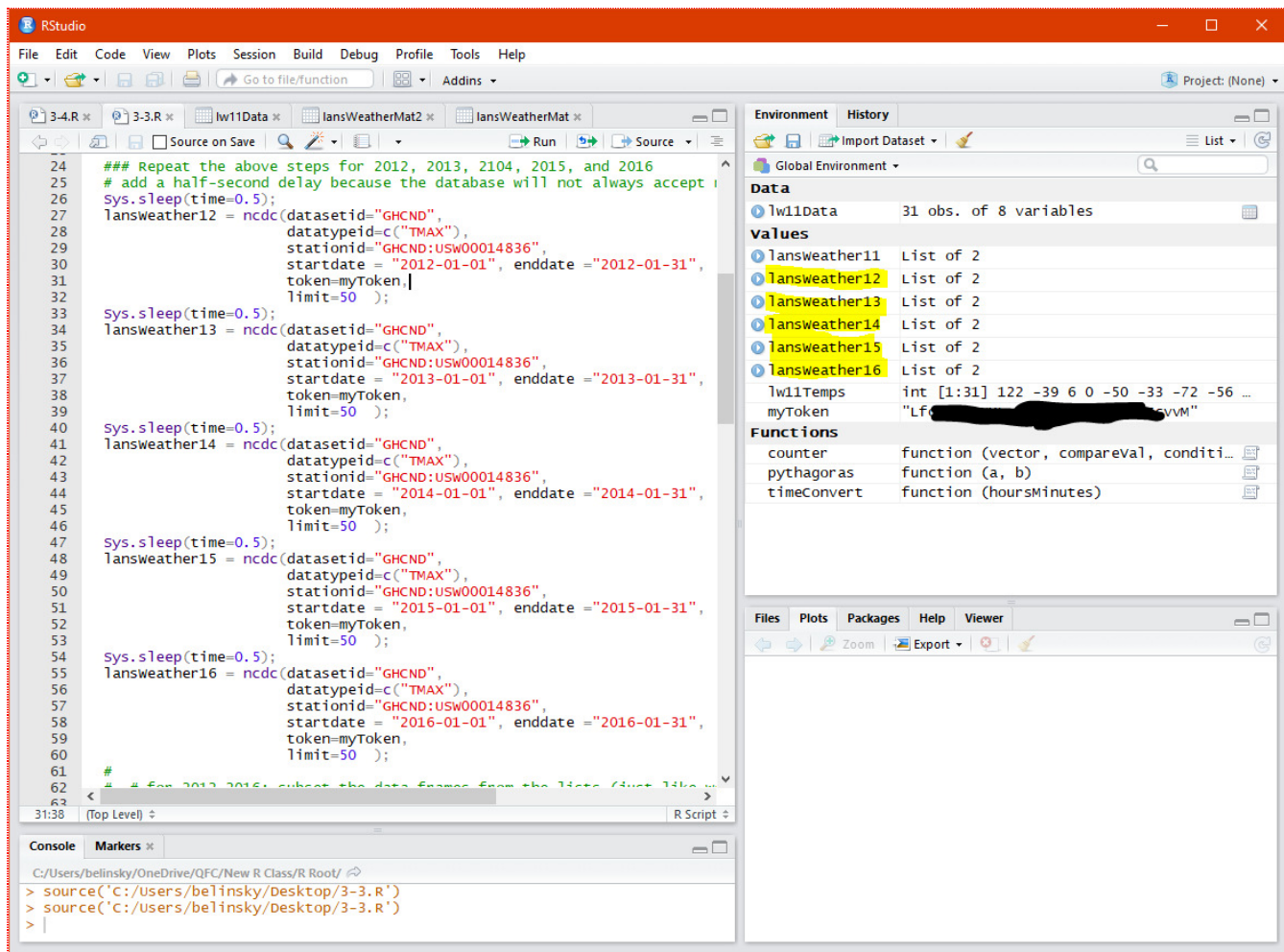


Fig 5: Getting temperature data from the other 5 years (2012-2016) with a half-second delay between requests

7 - Forming the matrix

A matrix is a two-dimensional structure (i.e., has rows and columns) where every value is of the same type. A matrix can also be thought of as a two-dimensional vector. In fact, one way to create a matrix is to break, or reshape, a vector up into rows and columns (i.e., make it two-dimensional).

We are going to do this by:

- Putting all the January temperatures collected from NOAA/NCDC into one large vector
- Reshape the vector into rows and columns

7.1 - Create one large vector

Let's put the six months of temperatures vector into one large vector:

```
1 lansTempsAll = c(lw11Temps, lw12Temps, lw13Temps,
2                 lw14Temps, lw15Temps, lw16Temps);
```

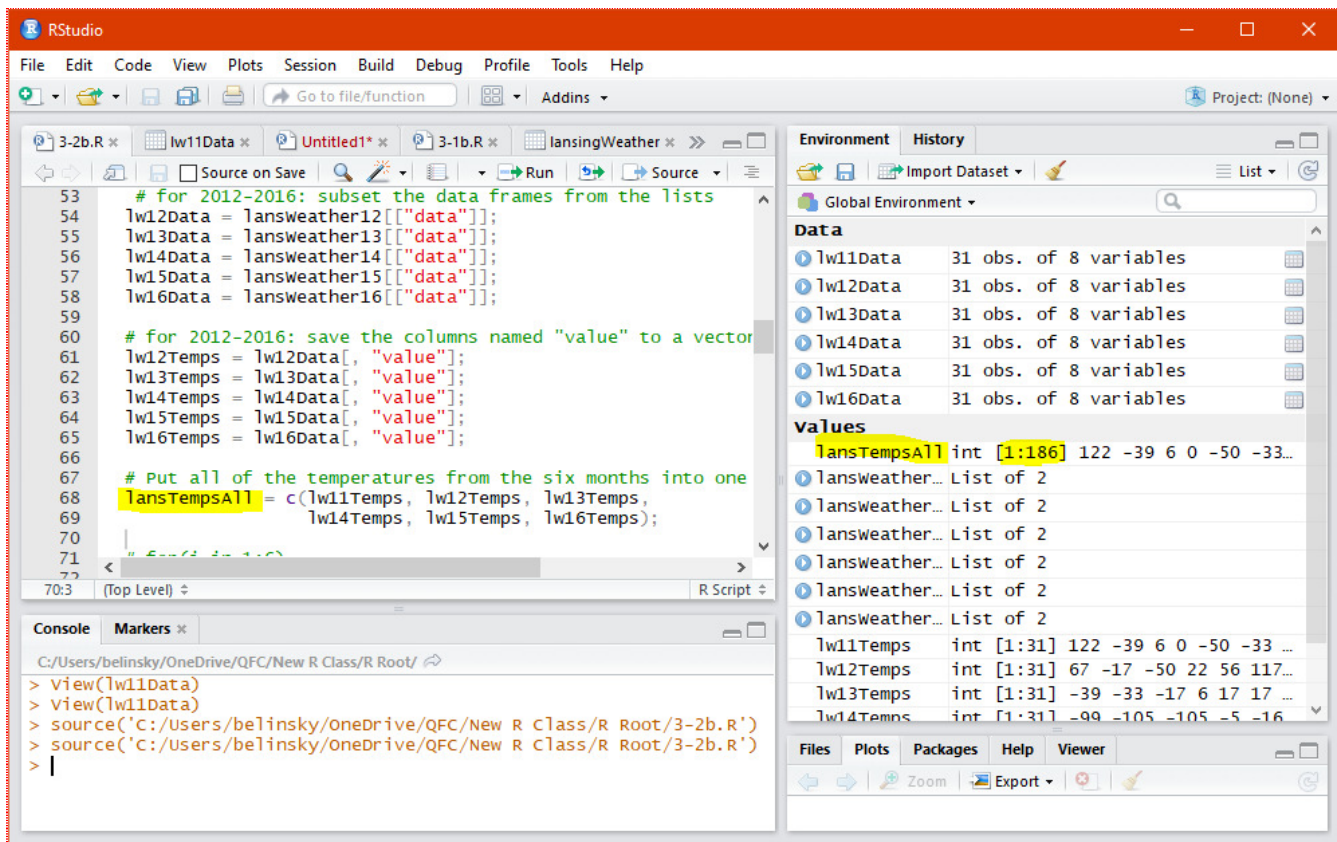



Fig 6: Vector that hold all 6 months of temperatures.

7.2 - Reshape the vector into rows and columns

Now you have a vector with $31 \times 6 = 186$ values. Our next step is to create a matrix from these 186 values. To do this, we call the function `matrix()`.

The parameters for `matrix()` are:

- **data**: the vector that holds the values you want to put in the matrix
- **nrow**: the number of rows
- **ncol**: the number of columns
- **byrow**: how the matrix is filled (**TRUE**: fill rows first, **FALSE**: fill columns first)

In this case, the rows will be the days (**nrow=31**) and the columns will be the years (**ncol=6**). But how do we fill the cells? We can either go across rows first (**byrow = TRUE**) or down columns first (**byrow = FALSE**). The temperature vector, `lansTempAll`, has all of January 2011 temperatures first, then all of January 2012 temperatures, etc. So, based on Fig 7, we want to go down columns first (**byrow = FALSE**).

	Jan 2011	Jan 2012	Jan2013	Jan2014	Jan 2015	Jan 2016
1	temp for Jan 1, 2011	temp for Jan 1, 2012	temp for Jan 1, 2013			
2	temp for Jan 2, 2011					
3	temp for Jan 3, 2011					
...						
30						
31						

Fig 7: Mock-up of the two-dimensional table that will hold all the January temperatures.

So the call to **matrix()** looks like this:

```
1 lansweatherMat = matrix(data=lansweatherAll, nrow=31, ncol=6, byrow = FALSE);
```

And the matrix looks like this:

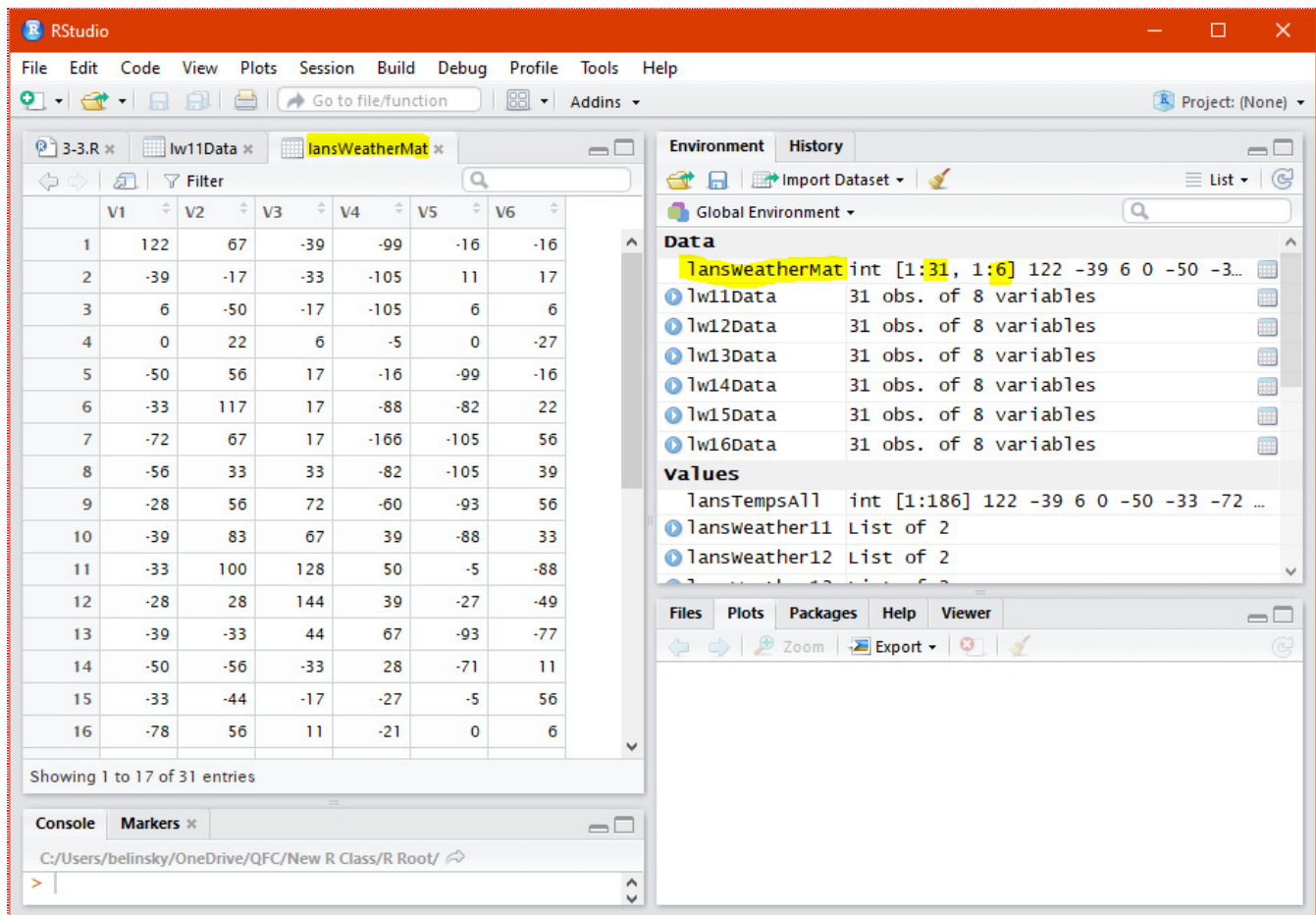


Fig 8: Matrix of 186 temperatures organized into 31 rows (days) and 6 columns (years).

Extension: Using `cbind()` and `rbind()` to create matrices ***Link does not exist***

7.3 - More about the parameters of `matrix()`

The number of values in a matrix equals the number of rows times the number of columns or:

`matrix.length = ncol * nrow`

If R is only given two of the values, R will calculate the third.

In other words, you can skip **`ncol`** or **`nrow`** and R will calculate the missing value for you assuming the number of values is known. This is useful if you don't know how much data you have (e.g., you are getting an unknown number of years worth of temperatures). R will just keep adding new columns as more years of data are added.

```
1 lansweatherMat2 = matrix(data=lansweatherAll, nrow=31, byrow = FALSE);
```

And, while there is not an intuitive use for it in this example (you are not adding days to January!), you can skip **`nrow`** and R will calculate it for you.

```
1 lansweatherMat3 = matrix(data=lansweatherAll, ncol=6, byrow = FALSE);
```

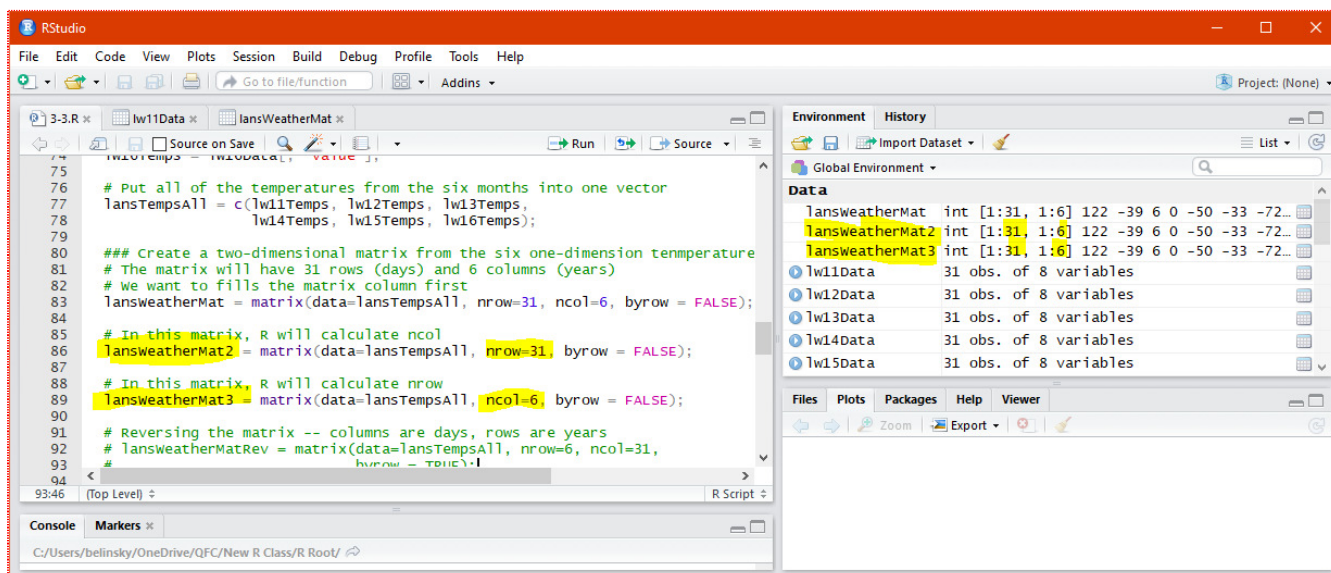


Fig 9: Having R calculate **nrow** or **ncol** -- the matrices are still the same.

Trap: When **nrow** * **ncol** does not equal the number of values.

7.4 - Reversing the matrix

In the previous example (Fig 9), the rows represent days and the columns represent years. We could easily reverse the matrix so that the **rows represent the years** and the **columns represent the days**. Visualizing this:

	1	2	3	...	30	31
Jan 2011	Jan 1, 2011	Jan 2, 2011	Jan 3, 2011			
Jan 2012	Jan 1, 2012					
Jan 2013	Jan 1, 2013					
Jan 2014						
Jan 2015						
Jan 2016						

Fig 10: Mock-up of the **REVERSED** two-dimensional table that will hold all the January temperatures

The temperature vector, **lansTempAll**, still has all of January 2011 temperatures first, then all of January 2012 temperatures, etc. So, based on Fig 10, we want to go **across rows first** (**byrow = TRUE**).

```
1 lansweatherMatRev = matrix(data=lansweatherAll, nrow=6, ncol=31, byrow = TRUE);
```


RStudio interface showing a matrix of temperature data. The matrix has 6 rows (years) and 31 columns (days). The values are in tenths of a Celsius degree.

	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	V14	V15	V16	V17	V18	V19	V20	V21	V22	V23	V24	V25	V26	V27	V28	V29	V30	V31
1	122	-39	6	0	-50	-33	-72	-56	-28	-39	-33	-28	-39	-50	-33	-78	-11	11	-56	-50	-100	-89	-106	-17	-17	-17	-28	-17	-33	-44	-78
2	67	-17	-50	22	56	117	67	33	56	83	100	28	-33	-56	-44	56	44	-44	-28	-83	-56	22	89	0	0	22	44	6	-17	50	117
3	-39	-33	-17	6	17	17	17	33	72	67	128	144	44	-33	-17	11	6	17	72	50	-94	-122	-83	-89	-67	-28	-11	61	139	133	-28
4	-99	-105	-105	-5	-16	-88	-166	-82	-60	39	50	39	67	28	-27	-21	-10	-60	-16	-10	-127	-116	-105	-71	-38	-38	-93	-132	-71	-21	-5
5	-16	11	6	0	-99	-82	-105	-105	-93	-88	-5	-27	-93	-71	-5	0	72	39	6	-5	-21	-10	0	28	6	-32	-27	-5	0	-49	6
6	-16	17	6	-27	-16	22	56	39	56	33	-88	-49	-77	11	56	6	-32	-93	-60	-55	-32	-16	-32	0	44	61	-16	11	-16	83	94

Fig 11: Transposing (i.e., reversing) the matrix: rows represent the 6 years and columns represent the 31 days

8 - Mathematics on the matrix

Let's go back to the original matrix with rows representing the days. Like last lesson, we have the issue that the temperatures are in tenths of a Celsius degree. To convert the temperatures to Fahrenheit we need to first divide the values by 10 and then convert from Celsius to Fahrenheit.

Since a matrix has all the same types of value, *we can apply these mathematical operations to the whole matrix* at once.

- 1 `lansweatherMat = lansweatherMat * 0.1;`
- 2 `lansweatherMat = (9/5) * lansweatherMat + 32;`

Note: This author is unsure why the values are rounded differently between the Environment Window and the matrix in the Main Window.

RStudio interface showing the matrix after conversion to Fahrenheit. The values are rounded to two decimal places. The Environment window shows the data structure.

	V1	V2	V3	V4	V5	V6
1	53.96	44.06	24.98	14.18	29.12	29.12
2	24.98	28.94	26.06	13.10	33.98	35.06
3	33.08	23.00	28.94	13.10	33.08	33.08
4	32.00	35.96	33.08	31.10	32.00	27.14
5	23.00	42.08	35.06	29.12	14.18	29.12
6	26.06	53.06	35.06	16.16	17.24	35.96
7	19.04	44.06	35.06	2.12	13.10	42.08
8	21.92	37.94	37.94	17.24	13.10	39.02
9	26.96	42.08	44.96	21.20	15.26	42.08
10	24.98	46.94	44.06	39.02	16.16	37.94
11	26.06	50.00	55.04	41.00	31.10	16.16
12	26.96	37.04	57.92	39.02	27.14	23.18
13	24.98	26.06	39.92	44.06	15.26	18.14

Environment window data structure:

- `lansweatherMat` num [1:31, 1:6] 54 25 33.1 32 23 ...
- `lansweatherMat2` int [1:31, 1:6] 122 -39 6 0 -50 -33 -72 -56..
- `lansweatherMat3` int [1:31, 1:6] 122 -39 6 0 -50 -33 -72 -56..
- `lansweatherMatR` int [1:6, 1:31] 122 67 -39 -99 -16 -16 -39 ...
- `lw11Data` 31 obs. of 8 variables
- `lw12Data` 31 obs. of 8 variables
- `lw13Data` 31 obs. of 8 variables
- `lw14Data` 31 obs. of 8 variables
- `lw15Data` 31 obs. of 8 variables

Fig 12: Performing a mathematical operation on the whole matrix -- all temperatures are now in Fahrenheit.

8.1 - Significant Digits

And, while we are at it, we only want two significant digits for each temperature (so, **53.96** becomes **54** and **9.14** become **9.1**). We can use the function **signif()** to set the significant digits.

signif() is similar to **round()** and has the parameters:

- **x**: values to apply the significant digits to
- **digits**: number of significant digits

```
1 | lansweatherMat = signif(x=lansweatherMat, digits=2);
```

Again, R adds **.0** to all integer numbers in the 4th column. This is because the 4th column was the only column that had values less than 10, so the second significant digit was a decimal. If any column has any decimal number, all numbers are given a decimal. The **.0** is not significant, numerically speaking.

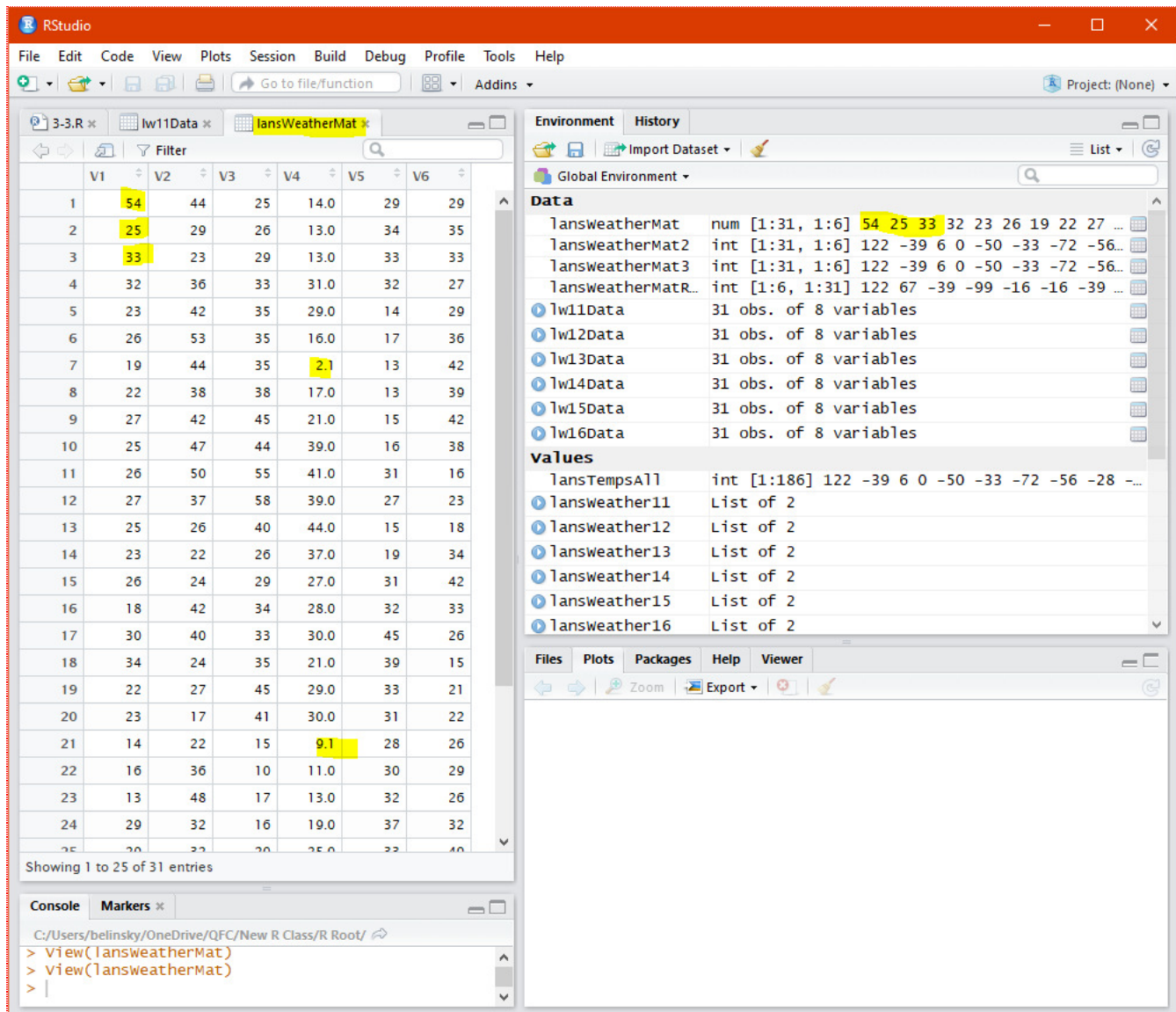


Fig 13: Setting the number of significant digits for the matrix.

9 - Statistics on the matrix

We can also apply statistical functions to the whole matrix. For instance, we can find the mean of all values in the matrix:

```
1 meanAllTemps = mean(lansWeatherMat);
```

We can find the mean of a specific column of the matrix, representing the mean of a certain year:

```
1 meanJan2013 = mean(lansWeatherMat[,3]);
```

Or find the mean of a specific row, representing all temperatures on a specific day in January:

```
1 meanJan17 = mean(lansWeatherMat[17,]);
```

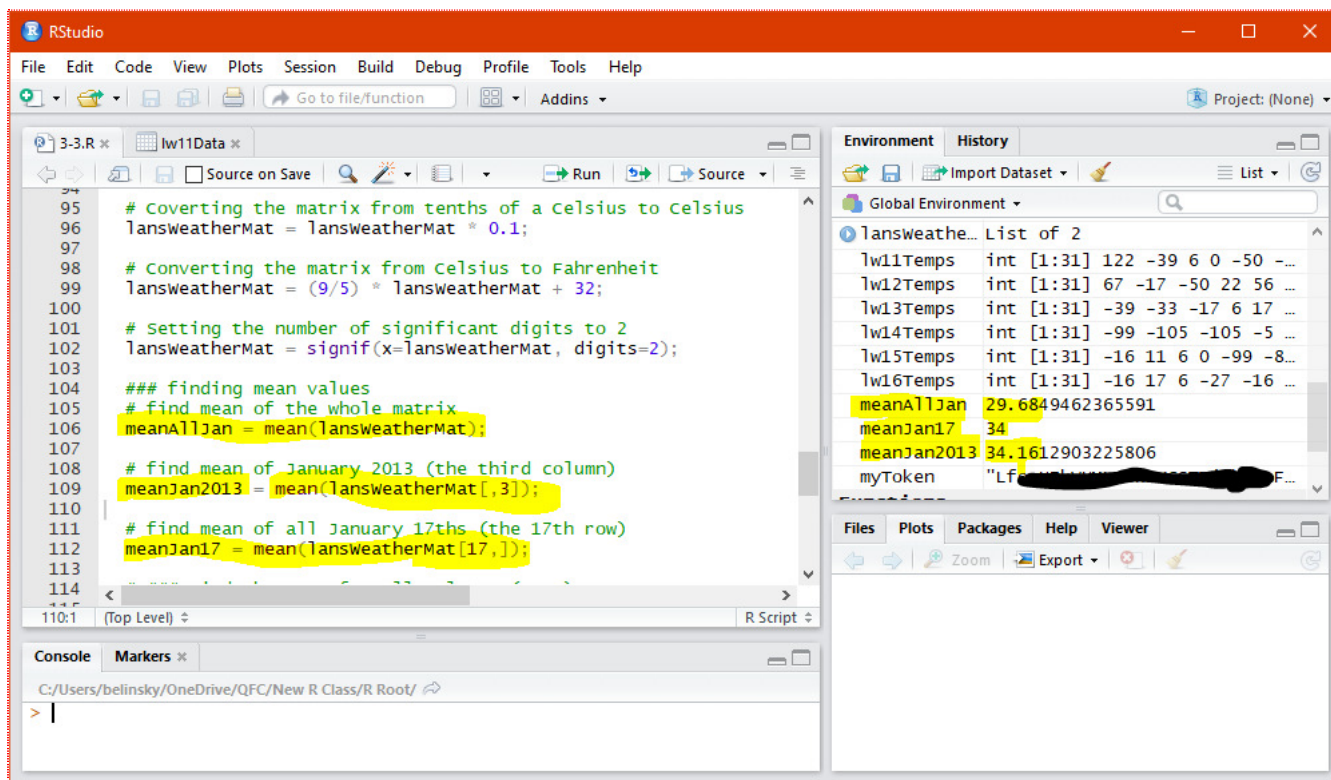


Fig 14: Finding means across a row, down a column, and for the whole matrix.

As a reminder, matrices, like a data frames, have two dimensions (rows and columns). So, we use two values to subset the matrix: a row and a column. If there is no subset value, then that means we are taking all values.

So: **[,3]** means take all rows in column 3
 and: **[17,]** means take all columns in row 17

9.1 - Performing a mean down each column

We are going to find the mean for all six years of January temperatures. The columns in **lansWeatherMatrix** each hold one January's worth of maximum temperatures. So we want to find six mean values, one for each column in the matrix.

We need to save the mean values to a vector. To do this, we create a variable, **yearlyMean**, as a vector that initially has no values.


```
1 yearlyMean = c(); # a vector that holds the 6 yearly mean temps
```

Then we will use a **for()** to iterate through the 6 columns (years). Remember that *i* will take on each of the six values in the sequence **1:6** in turn as the loop iterates.

```
1 for(i in 1:6)
2 {
3 }
```

In the codeblock attached to the **for()**, we get the mean value of the temperatures in the indexed column: **mean(JanTempMatrix[,i])**, where *i* will take on the values **1-6**.

Note: **JanTempMatrix[,i]** means get the value from every row in column *i*.

We then assign the mean value to the **yearlyMean** vector: **yearlyMean[i]**

```
1 yearlyMean = c();
2
3 for(i in 1:6)
4 {
5   # get the mean of all values in column i and save it to monthlyMean[i]
6   yearlyMean[i] = mean(lansWeatherMatrix[,i]);
7 }
```

Once again, R displays numbers with different numbers of decimal places in different windows -- but the numbers are the same.

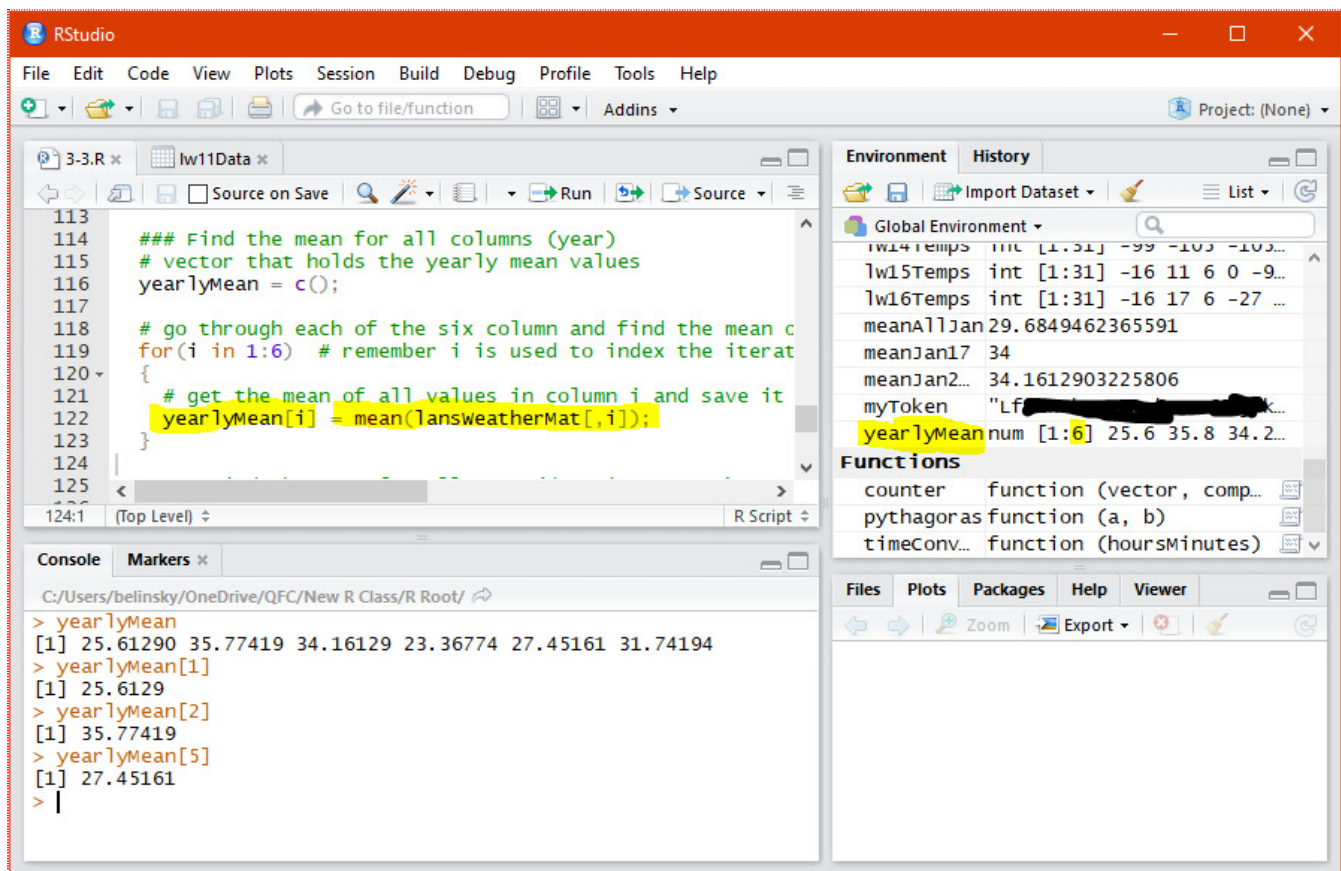


Fig 15: Getting the mean temperature for each month using a **for()**

9.2 - Performing a mean across all rows (days)

We can do the same thing across all rows to find the average temperature for each day in January across the six years.

In this case, we need to iterate the **for()** 31 times (**1:31**), and change the subset operation from

```
1 | mean(lansweatherMatrix[,i]); # get all row values from column i
   |
   | to
1 | mean(lansweatherMatrix[i,]); # get all column values from row i
```

```
1 | dailyMean = c();
2 |
3 | for(i in 1:31)
4 | {
5 |   # get the mean of all values in row i and save it to dailyMean[i]
6 |   dailyMean[i] = mean(lansweatherMatrix[i,]);
7 | }
```

We can use the Console Window to display any mean value. For example typing "dailyMean[13]" in the Console Window will output the **13th** value in **dailyMean**.

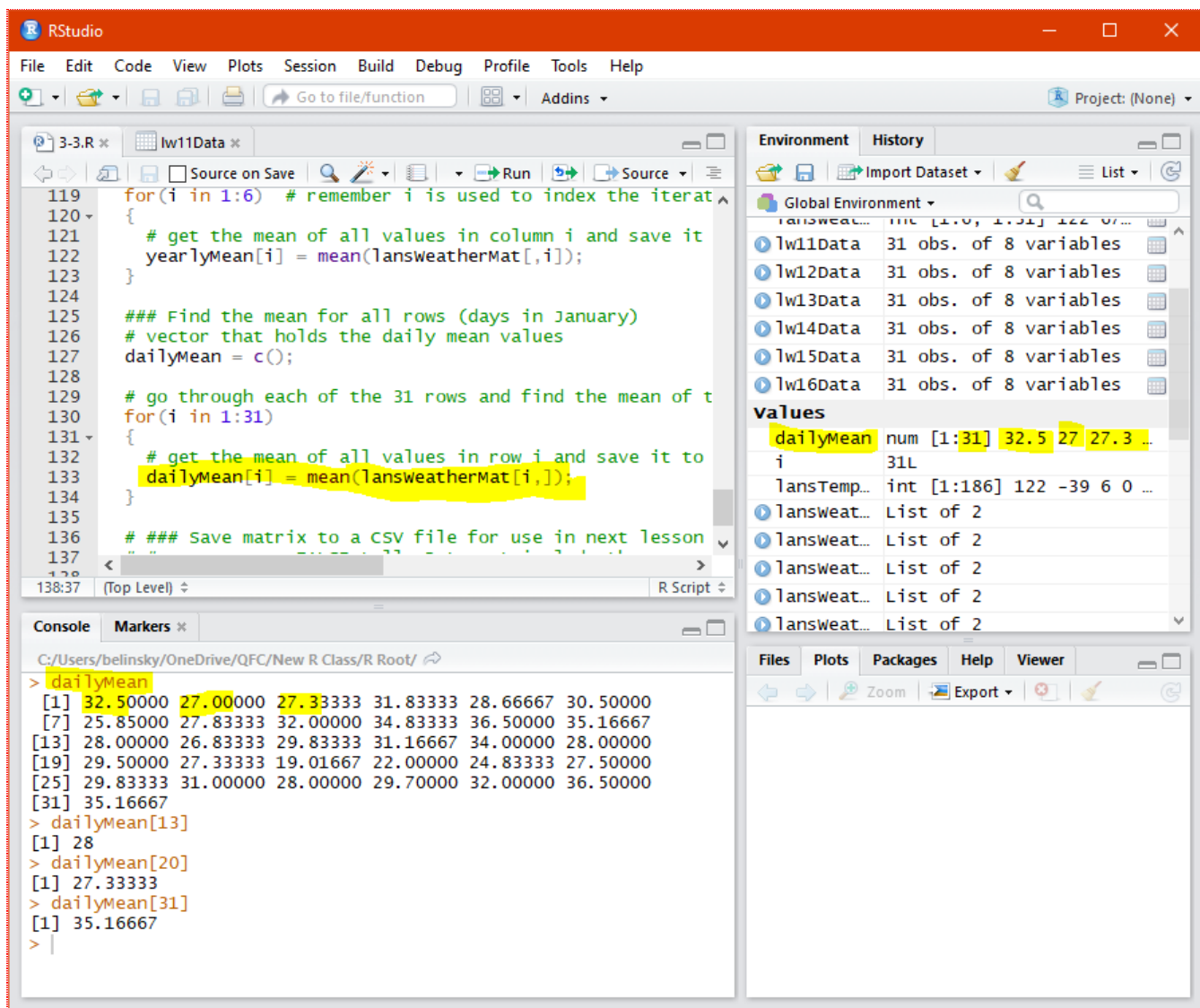


Fig 16: Getting the mean temperature for each day in January using `for()`

10 - Saving the matrix

We are going to save ***lansWeatherMatrix*** to a CSV file so we can use it in the next lesson and avoid the hassle of calling NOAA/NCDC and reformatting the data.

```
1 write.csv(x=lansWeatherMat, file = "data/LansingJanTemps.csv", row.names = FALSE);
```

Now you should have a file called ***LansingJanTemps.csv*** in the ***data*** folder inside your ***R Root*** directory.

11 - Application

1) Get the wind speed from the NOAA/NCDC database for January, February, March, and April 2016 in Lansing, MI.

- wind speed is ***datatypeid: AWND***
- query each month separately

2) Save the wind speeds into 4 vectors, each holding one months worth of wind values

- 3) Make all vectors equal in length (31) by added NA values to the end of February and April.
- 4) Combine the four vectors into one vector.
- 5) Create a matrix from the combined vector that has the months as columns and the days as rows:

	January	February	March	April
1				
2				
...				
30		NA		
31		NA		NA

- 6) Find the mean, standard deviation, and variance for the wind speed for February and March
 - you will need to use the parameter **na.rm** to deal with the **NA** values.
- 7) Find the month that has the highest mean wind speed.
- 8) Find the day that has the highest mean wind speed.

12 - Extension: Using `cbind()` and `rbind()` to create matrices

After the six calls to the NOAA/NCDC database, we extracted six vectors, each vector represented one January worth of maximum temperatures -- from 2011 to 2016.

In this lesson, we names those vector variables: ***lw11Temps***, ***lw12Temps***, ***lw13Temps***, ***lw14Temps***, ***lw15Temps***, and ***lw16Temps***

For the data in this lesson, there is a quick way to turn these six vector into a matrix and that is using ***cbind()*** or ***rbind()***. ***cbind()*** binds vectors together into a matrix placing each vector in a new column. ***rbind()*** binds vectors together into a matrix placing each vector in a new row. Effectively, ***rbind()*** created the transpose matrix of ***cbind()***.

```

1  cbindMatrix = cbind(lw11Temps, lw12Temps, lw13Temps,
2                        lw14Temps, lw15Temps, lw16Temps);
3
4  rbindMatrix = rbind(lw11Temps, lw12Temps, lw13Temps,
5                      lw14Temps, lw15Temps, lw16Temps);

```

	lw11Temps	lw12Temps	lw13Temps	lw14Temps	lw15Temps	lw16Temps
1	122	67	-39	-99	-16	-16
2	-39	-17	-33	-105	11	17
3	6	-50	-17	-105	6	6
4	0	22	6	-5	0	-27
5	-50	56	17	-16	-99	-16
6	-33	117	17	-88	-82	22
7	-72	67	17	-166	-105	56
8	-56	33	33	-82	-105	39
9	-28	56	72	-60	-93	56
10	-39	83	67	39	-88	33
11	-33	100	128	50	-5	-88
12	-28	28	144	39	-27	-49
13	-39	-33	44	67	-93	-77
14	-50	-56	-33	28	-71	11
15	-33	-44	-17	-27	-5	56

Showing 1 to 15 of 31 entries

Fig 17: `cbind()` makes each vector a column in the new matrix

3-3.r

3-4.r

3-5.r

3-6.r

3-3.R

cbindMatrix

rbindMatrix

lw11Data

Filter

lw11Temps	122	-39	6	0	-50	-33	-72	-56	-28	-39	-33	-28	-39	-50	-33	-78	-11	11	-56	-50
lw12Temps	67	-17	-50	22	56	117	67	33	56	83	100	28	-33	-56	-44	56	44	-44	-28	-83
lw13Temps	-39	-33	-17	6	17	17	17	33	72	67	128	144	44	-33	-17	11	6	17	72	50
lw14Temps	-99	-105	-105	-5	-16	-88	-166	-82	-60	39	50	39	67	28	-27	-21	-10	-60	-16	-10
lw15Temps	-16	11	6	0	-99	-82	-105	-105	-93	-88	-5	-27	-93	-71	-5	0	72	39	6	-5
lw16Temps	-16	17	6	-27	-16	22	56	39	56	33	-88	-49	-77	11	56	6	-32	-93	-60	-55

Fig 18: `rbind()` makes each vector a row in the new matrix

13 - Trap: `nrow * ncol` does not equal number of values