

Plot and Subset Precipitation Data in R - 2013

Colorado Floods

In this lesson, we will learn how to import a larger dataset, and test our skills cleaning and plotting the data.

Learning Objectives

After completing this tutorial, you will be able to:

- Import a text file into R.
- Plot quantitative time series data using `ggplot`
- Ensure that NoData values do not interfere with quantitative analysis by setting them to `NA` in R.
- Use the `na.rm` argument when performing math with large datasets.
- Subset data using the `dplyr filter()` function
- Use `dplyr` pipes to filter data in R.

What you need

You need R and RStudio to complete this tutorial. Also you should have an `earth-analytics` directory setup on your computer with a `/data` directory with it.

- How to Setup R / RStudio
- Setup your working directory
- Intro to the R & RStudio Interface

R Libraries to Install:

- **ggplot2:** `install.packages("ggplot2")`
- **dplyr:** `install.packages("dplyr")`

Download Week 2 Data{data-proofer-ignore=} .btn }

Important - Data Organization

Before you begin this lesson, be sure that you've downloaded the dataset above. You will need to UNZIP the zip file. When you do this, be sure that your directory looks like the image below: note that all of the data are within the `week2` directory. They are not nested within another directory. You may have to copy and paste your files to make this look right.

Your `week2` file directory should look like the one above. Note that the data directly under the `week-2` folder.

Get started with time series data

Let's get started by loading the `ggplot2` and `dplyr` libraries. Also, let's set our working directory. Finally, set `stringsAsFactors` to `FALSE` globally as shown below.

```
# set your working directory to the earth-analytics directory
# setwd("working-dir-path-here")

# load packages
```

```
library(ggplot2) # efficient plotting
library(dplyr) # efficient data manipulation

# set strings as factors to false for everything
options(stringsAsFactors = FALSE)
```

Import precipitation time series

We will use a precipitation dataset derived from data accessed through the National Centers for Environmental Information (formerly National Climate Data Center) Cooperative Observer Network (COOP) station 050843 in Boulder, CO. The data time span is: 1 January 2003 through 31 December 2013.

We can use `read.csv()` to import the `.csv` file.

```
# download the data
# download.file(url = "https://ndownloader.figshare.com/files/7283285",
#               destfile = "data/week2/805325-precip-dailysum_2003-2013.csv")

# import the data
boulder_daily_precip <- read.csv("data/week2/precipitation/805325-precip-dailysum-2003-2013.csv",
                                header = TRUE)

# view first 6 lines of the data
head(boulder_daily_precip)
##      DATE DAILY_PRECIP   STATION  STATION_NAME ELEVATION LATITUDE
## 1 1/1/03         0.00 COOP:050843 BOULDER 2 CO US    1650.5 40.03389
## 2 1/5/03        999.99 COOP:050843 BOULDER 2 CO US    1650.5 40.03389
## 3 2/1/03         0.00 COOP:050843 BOULDER 2 CO US    1650.5 40.03389
## 4 2/2/03        999.99 COOP:050843 BOULDER 2 CO US    1650.5 40.03389
## 5 2/3/03         0.40 COOP:050843 BOULDER 2 CO US    1650.5 40.03389
## 6 2/5/03         0.20 COOP:050843 BOULDER 2 CO US    1650.5 40.03389
##  LONGITUDE YEAR JULIAN
## 1 -105.2811 2003      1
## 2 -105.2811 2003      5
## 3 -105.2811 2003     32
## 4 -105.2811 2003     33
## 5 -105.2811 2003     34
## 6 -105.2811 2003     36

# view structure of data
str(boulder_daily_precip)
## 'data.frame':    792 obs. of  9 variables:
##  $ DATE       : chr  "1/1/03" "1/5/03" "2/1/03" "2/2/03" ...
##  $ DAILY_PRECIP: num   0e+00 1e+03 0e+00 1e+03 4e-01 ...
##  $ STATION     : chr  "COOP:050843" "COOP:050843" "COOP:050843" "COOP:050843" ...
##  $ STATION_NAME: chr  "BOULDER 2 CO US" "BOULDER 2 CO US" "BOULDER 2 CO US" "BOULDER 2 CO US" ...
##  $ ELEVATION   : num   1650 1650 1650 1650 1650 ...
##  $ LATITUDE    : num    40 40 40 40 40 ...
##  $ LONGITUDE   : num  -105 -105 -105 -105 -105 ...
##  $ YEAR        : int   2003 2003 2003 2003 2003 2003 2003 2003 2003 2003 ...
##  $ JULIAN      : int    1 5 32 33 34 36 37 38 41 49 ...
```

```
# are there any unusual / No data values?
summary(boulder_daily_precip$DAILY_PRECIP)
##      Min.   1st Qu.    Median      Mean   3rd Qu.      Max.
##    0.000    0.100    0.100    5.297    0.300   1000.000
max(boulder_daily_precip$DAILY_PRECIP)
## [1] 999.99
```

About the Data

Viewing the structure of these data, we can see that different types of data are included in this file.

- **STATION** and **STATION_NAME**: Identification of the COOP station.
- **ELEVATION**, **LATITUDE** and **LONGITUDE**: The spatial location of the station.
- **DATE**: The date when the data were collected in the format: YYYYMMDD. Notice that DATE is currently class `chr`, meaning the data is interpreted as a character class and not as a date.
- **DAILY_PRECIP**: The total precipitation in inches. Important: the metadata notes that the value 999.99 indicates missing data. Also important, hours with no precipitation are not recorded.
- **YEAR**: the year the data were collected
- **JULIAN**: the JULIAN DAY the data were collected.

Additional information about the data, known as metadata, is available in the `PRECIP_HLY_documentation.pdf`. The metadata tell us that the noData value for these data is 999.99. IMPORTANT: we have modified these data a bit for ease of teaching and learning. Specifically, we've aggregated the data to represent daily sum values and added some noData values to ensure you learn how to clean them!

You can download the original complete data subset with additional documentation [here](#).

Challenge

Using everything you've learned in the previous lessons:

- Import the dataset: `data/week2/precipitation/805325-precip-dailysum-2003-2013.csv`
- Clean the data by assigning noData values to NA
- Make sure the date column is a date class
- When you are done, plot it using `ggplot()`.
- Be sure to include a TITLE, and label the X and Y axes.
- Change the color of the plotted points

Some notes to help you along:

- Date: be sure to take of of the date format when you import the data.
- NoData Values: We know that the no data value = 999.99. We can account for this when we read in the data. Remember how?

Your final plot should look something like the plot below.

****Data Tip:****For a more thorough review of date/time classes, see the NEON tutorial *Dealing With Dates & Times in R - as.Date, POSIXct, POSIXlt*. {`: notice`}

Optional challenge

Take a close look at the plot.

- What does each point represent?

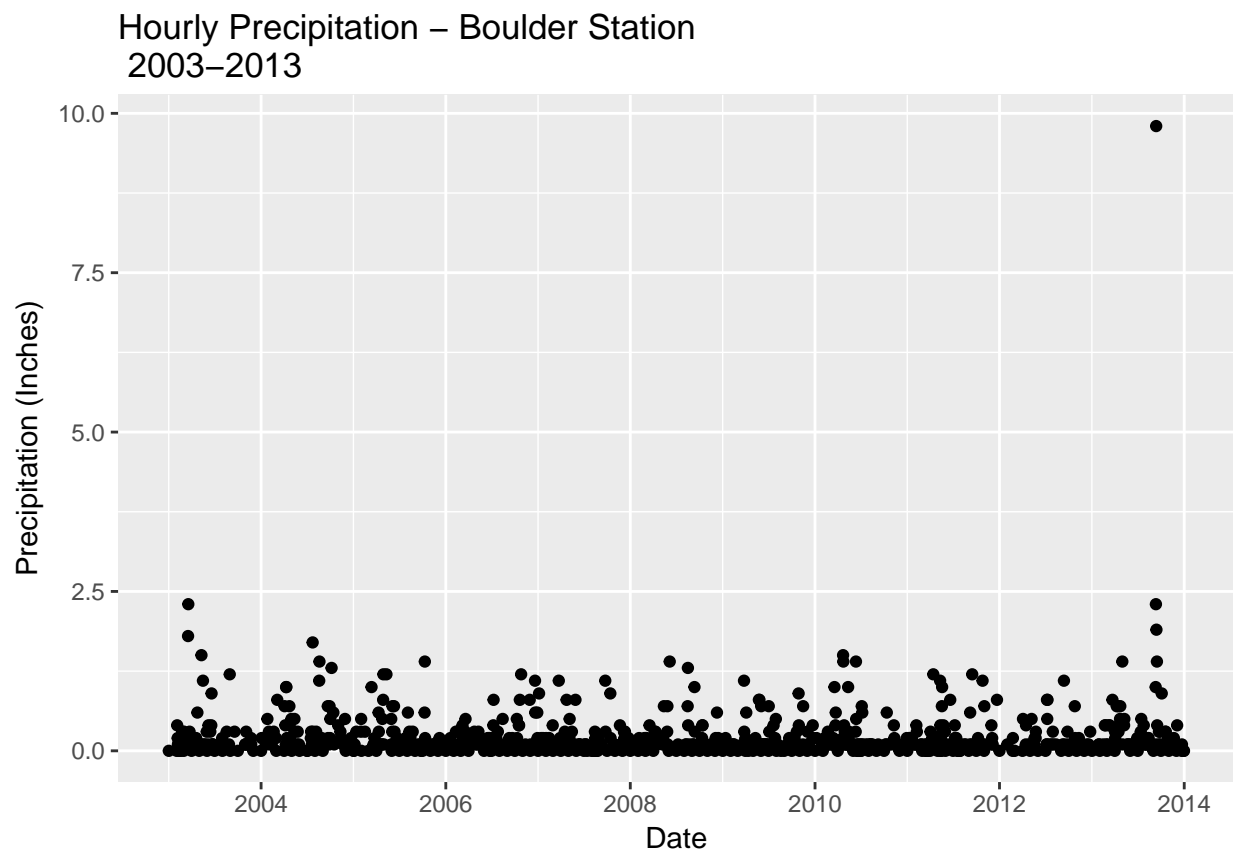


Figure 1: precip plot w fixed dates

- Use the `min()` and `max()` functions to determine the minimum and maximum precipitation values for the 10 year span?

Subset the Data

If we wanted to zoom in and look at some data over a smaller time period, we can subset it. Let create a subset of data for the time period around the flood between 15 August to 15 October 2013. We will use the `filter()` function in the `dplyr` package to do this.

Introduction to the pipe `%>%`

Pipes let you take the output of one function and send it directly to the next, which is useful when you need to do many things to the same data set. Pipes in R look like `%>%` and are made available via the `magrittr` package, installed automatically with `dplyr`.

```
# subset 2 months around flood
precip_boulder_AugOct <- boulder_daily_precip %>%
  filter(DATE >= as.Date('2013-08-15') & DATE <= as.Date('2013-10-15'))
```

In the code above, we use the pipe to send the `boulder_daily_precip` data through a filter step. In that filter step, we filter out only the rows within the date range that we specified. Since `%>%` takes the object on its left and passes it as the first argument to the function on its right, we don't need to explicitly include it as an argument to the `filter()` function.

```
# check the first & last dates
min(precip_boulder_AugOct$DATE)
## [1] "2013-08-21"
max(precip_boulder_AugOct$DATE)
## [1] "2013-10-11"

# create new plot
precPlot_flood2 <- ggplot(data=precip_boulder_AugOct, aes(DATE,DAILY_PRECIP)) +
  geom_bar(stat="identity") +
  xlab("Date") + ylab("Precipitation (inches)") +
  ggtitle("Daily Total Precipitation Aug - Oct 2013 for Boulder Creek")

precPlot_flood2
```

Challenge

Create a subset from the same dates in 2012 to compare to the 2013 plot. Use the `ylim()` argument to ensure the y axis range is the SAME as the previous plot - from 0 to 10“.

How different was the rainfall in 2012?

HINT: type `?lims` in the console to see how the `xlim` and `ylim` arguments work.

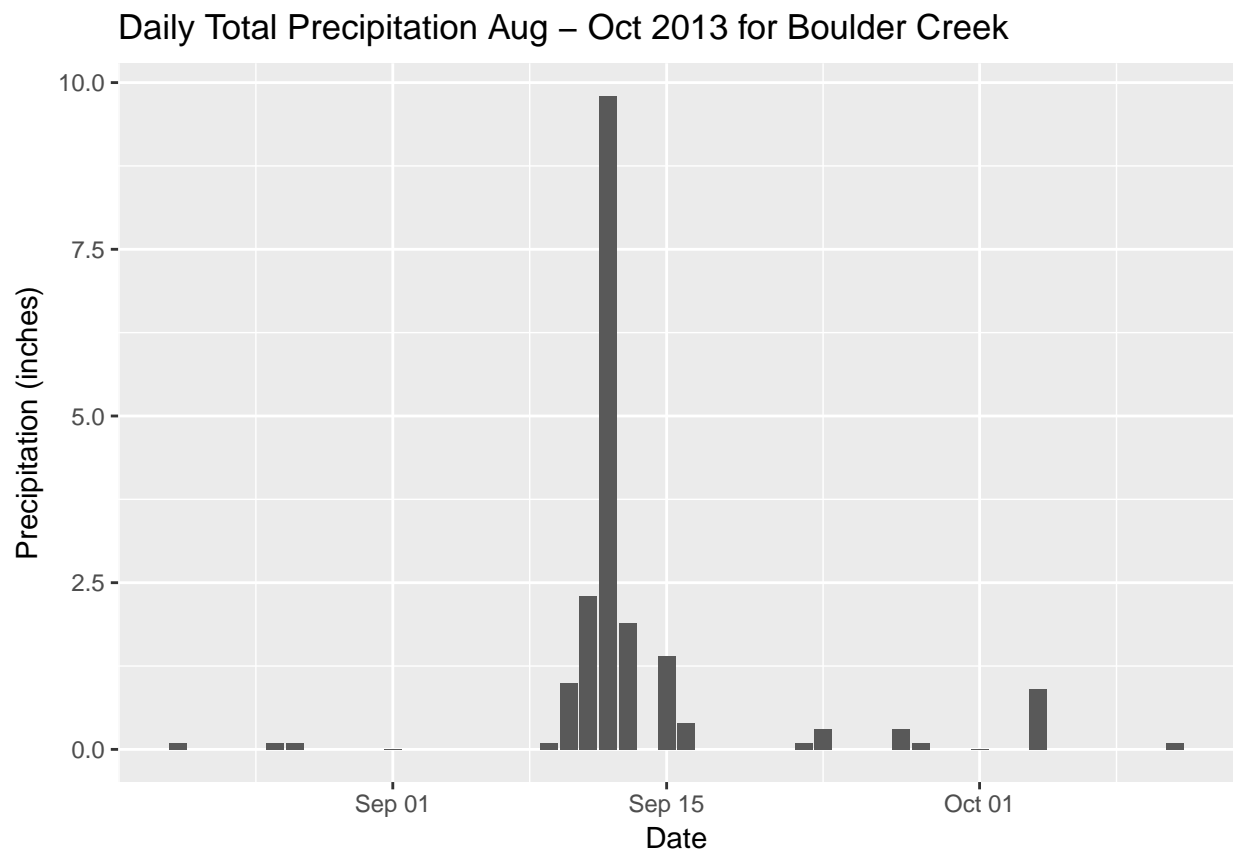


Figure 2: precip plot subset

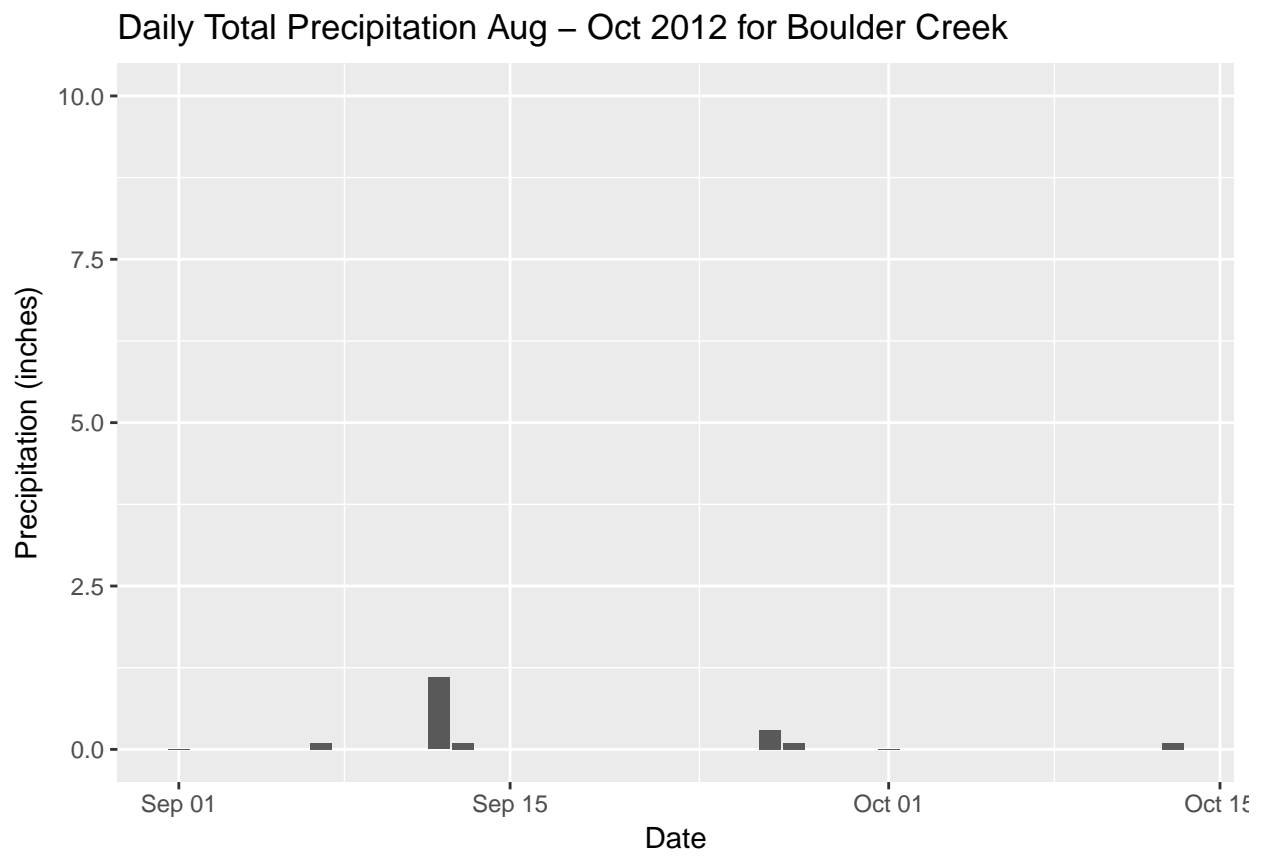


Figure 3: precip plot subset 2