

# Guide 2: Cleaning and Pre-Processing Weather Station Data

Marc Los Huertos

February 1, 2025 (ver. 0.95)

## 1 Introduction

### 1.1 Goals

The goal of this guide is to provide a step-by-step process for cleaning and pre-processing weather station data. The data is from the Global Historical Climatology Network (GHCN) Daily dataset. The data is available from the National Oceanic and Atmospheric Administration (NOAA) and is available from the National Centers for Environmental Information (NCEI) at <https://www.ncei.noaa.gov/products/land-based-station/global-historical-climatology-network-daily>.

### 1.2 Background

#### 1.2.1 What is GHCN-Daily (GNCN-d)?

The new dataset is referred to as GHCN-Daily, version 3. The new dataset has a number of improvements over the older dataset, including a more comprehensive and robust quality assurance process. The new dataset also includes more station records and is updated more frequently than the older dataset.

The new dataset is also available in a more user-friendly format, including a set of pre-packaged data files that are available for download from the NCEI website. However, developing instructions to interact with their website is way to time consuming! So, we'll use the R file transfer protocol, via https, to download the data.

In 2011, there was a major update to the daily updates to the GHCN dataset: [IMMEDIATE – Changes to COOP Daily Form Access](#). It's a useful history.

### 1.3 Approach

We need to address several things that might get in the way of our analysis:

1. Read csv files into R (from station.csv files)
2. Fix date format (convert to POSIXct, adding month and year as variables)

3. Convert units (VALUE: PRCP = 0.1 mm, TMAX/TMIN = 0.1 C)
4. Evaluate missing data (Coverage.fun)
5. Evalutae for Outliers (QACQ.fun)
6. Other stuff??

## 2 Cleaning and Pre-Processing Functions

### 2.1 Before Starting the Process

Before you begin, make sure you have the stations to read into R. Got to the “Files” tab in RStudio and make sure you see the station csv files in your data folder. If not, please Slack Marc and mentors, so we can troubleshoot the issue!

### 2.2 R Code with Custom Functions

From the Canvas page, go to the Guide2functions.R file and download the file to your computer. Then upload the file to Rstudio directory you are using for the project.

Open the file in Rstudio and run the code, using the “source”. button near the top of the editor window.

Run the **Guide2functions.R** code and the functions will be loaded into your environment automatically. Be sure the function has been updated to 2025-02-01!

### 2.3 Function Descriptions and Use

The following custom functions are used to read the weather stations data into R, clean and pre-process data so we can analyze the data with the next guide.

### 2.4 Reading csv and Checking dataframe objects

**Read all csv files into R** The data are now ready to be read into R environment, we will read them in using the following function:

**Function:** **ReadStations2.fun()**

Use this function read the stations in the R Environment.

Example of how to use the function:

```
datafolder = "/home/mwl04747/RTricks/05_Regional_Climate_Trends/Data/SP25/"
ReadStations2.fun(datafolder)
```

**List the objects in the R environment** The `ls()` function will list the objects in the R environment. Use this in the console.

```
ls()

## [1] "coverage.fun"      "datafolder"      "df_exists"
## [4] "df_head"           "df_names"        "fixDates.fun"
## [7] "fixValues.fun"     "MonthlyAnomalies.fun" "MonthlyNormals.fun"
## [10] "MonthlyValues.fun" "QAQC.fun"        "ReadStations2.fun"
## [13] "SaveCleanUp.fun"   "sortStations.fun" "USC00040693"
## [16] "USC00041614"       "USC00042294"     "USC00043157"
## [19] "USC00044259"       "USC00044412"     "USC00046074"
## [22] "USC00046136"       "USC00046719"     "USC00046826"
## [25] "USC00047821"       "USC00047902"     "USC00048351"
## [28] "USC00049866"       "USW00023271"
```

What do you see? Your list of objects should include several things like mine: some weather stations, You should see objects named after the stations that have been read into R. In addition, you should see several functions (.fun) and other miscellaneous objects. If so, you have been making progress. If you don't see something parallel to my list, let me or mentors know so we can trouble shoot, or look at section 4.

**Check the structure of the dataframes** Using `str()`, make sure the datasets look right!

For example:

```
## 'data.frame': 359900 obs. of 8 variables:
## $ ID : chr "USC00042294" "USC00042294" "USC00042294" "USC00042294" ...
## $ DATE : int 18930101 18930102 18930103 18930104 18930105 18930106 18930107 18930108 ...
## $ ELEMENT : chr "TMAX" "TMAX" "TMAX" "TMAX" ...
## $ VALUE : int 283 100 89 89 83 83 83 78 89 89 ...
## $ M.FLAG : chr "" "" "" "" ...
## $ Q.FLAG : chr "" "" "" "" ...
## $ S.FLAG : chr "6" "6" "6" "6" ...
## $ OBS.TIME: int NA NA NA NA NA NA NA NA NA NA ...
```

What to look for? Is the object a data.frame? Does it have the right number of observations? Are the expected variable names present? Are the data types correct? Are values in each variable reasonable?

**Sorting Stations by Number of Observations** In theory, the stations with the greatest number of observations are the ones we will want to evaluate. However, we may also want to evaluate stations to cover a geographic range.

THIS IS A BRAND NEW IDEA FROM WEDNESDAY’S LAB, AND WILL WORK ON THIS FOR NEXT SEMESTER. Instead look at the number of observation in the R environment to select 3 stations with the highest number of observations.

**Function:** `sortstations.fun()`

Example of how to use the function:

```
# sortstations.fun() Nor working yet.
```

## 2.5 Clean Data

Next we’ll “clean” the dataset by fixing the date format and preparing it for the analysis stages. I suggest you select two stations with the highest number of observations to work with. We will know if this is a problem soon enough and if it is, we can use one of the 13 remaining stations, if your state/territory has that many to choose from.

**Function to Fix Dates** This function converts date values to a format that R can understand as a date, so we can, among other things, create montly means.

**Function:** `fixDates.fun()`

Example of how to use the function (Note: the new dataframe name change with an a following the station.):

```
USC00042294a <- fixDates.fun(USC00042294)
USC00040693a <- fixDates.fun(USC00040693)
```

**Evaluation Data Coverage** We need to know how much data we have for each station. This is important for the next steps in the process.

**Function:** `coverage.fun()`

Example of how to use the function:

```
coverage.fun(USC00042294a)
coverage.fun(USC00040693a)
```

In general, we want something like 95% coverage for the period of record. If you don’t have that, please let us know and we’ll help you get additional stations with better coverage!

See Section ?? for more information on how to evaluate the coverage of the data.

**Function to Fix Values (by order of magnitude)** According to the NOAA website, the csv.gz files have five core elements include the following units, which we will convert:

**PRCP** = Precipitation (tenths of mm) → mm

**SNOW** = Snowfall (mm) → cm

**SNWD** = Snow depth (mm) → cm

**TMAX** = Maximum temperature (tenths of degrees C) → degrees C

**TMIN** = Minimum temperature (tenths of degrees C) → degrees C

**Function:** `fixValues.fun()`

Example of how to use the function (Note: the new dataframe name change):

```
USC00042294b <- fixValues.fun(USC00042294a)
USC00040693b <- fixValues.fun(USC00040693a)
```

**Checking for Outliers** NOAA website conducts a very rudimentary data quality check, but every year, we find stations with wacky numbers. Hopefully this function will find them. But if you do have some, let Marc and mentors know so we can figure out how to address them!

Here are some stuff we'll look for:

**Extreme Values** Plot values with time, is the scale crazy with just a few observations at the extreme?

**Sudden Shift** ???

**Note QA/QC Flags** TBD

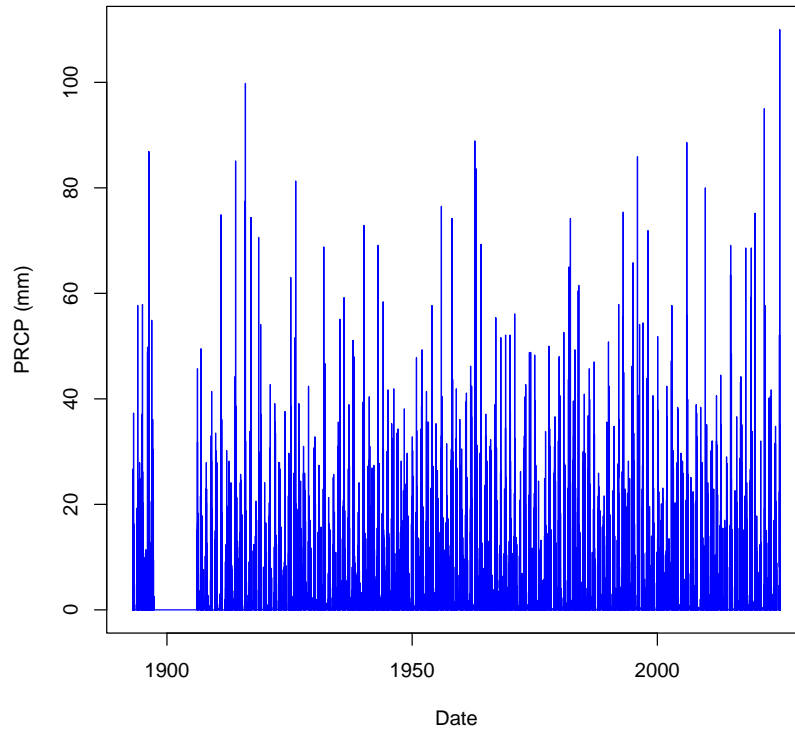
Yes, I am still working on this. QA/QC is a big deal and I want to make I can explain it AND make sure the code works.

**Function:** `QAQC.fun()`

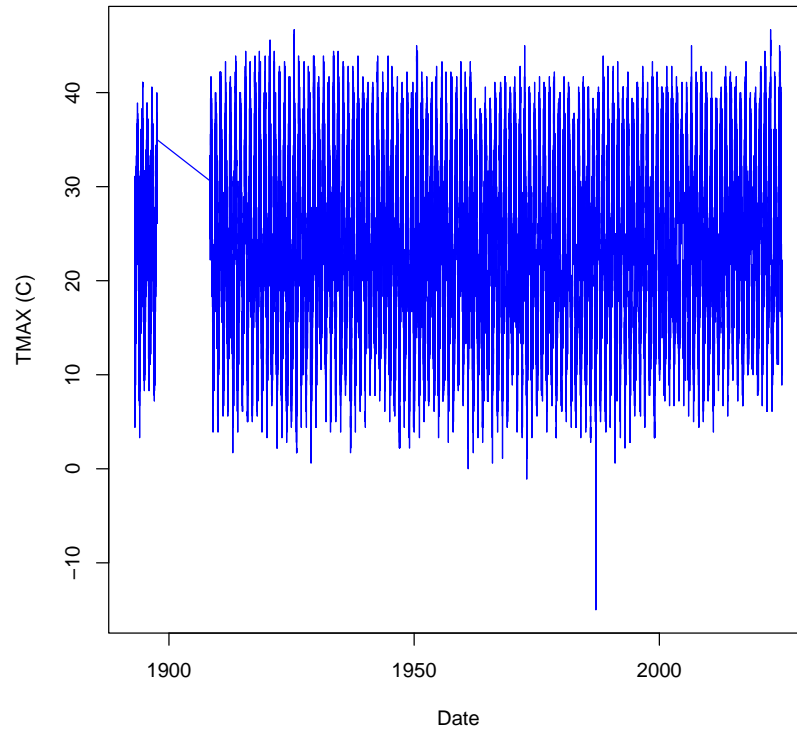
Example of how to use the function:

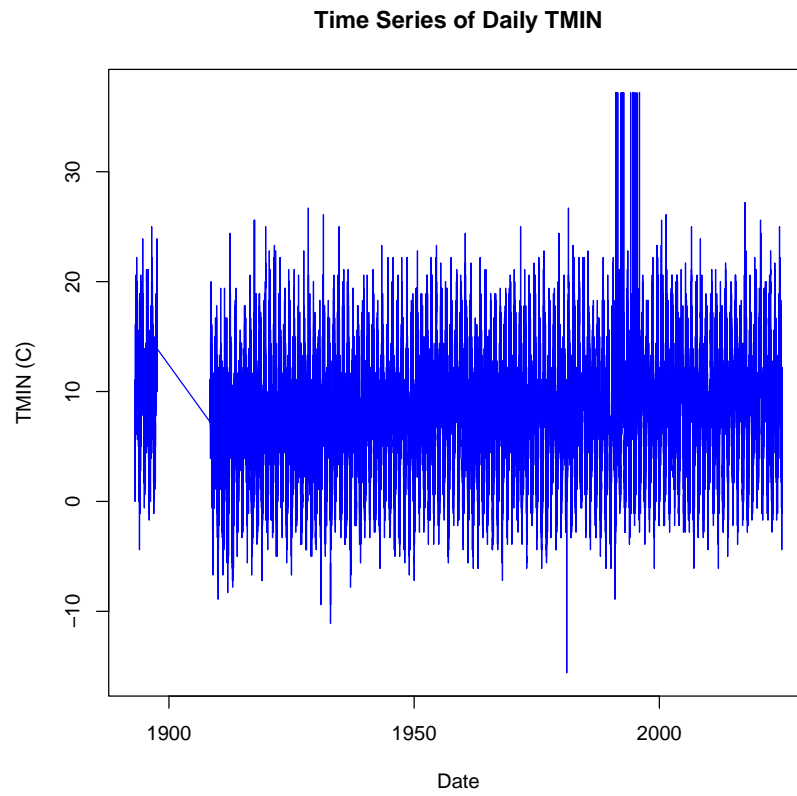
```
QAQC.fun(USC00042294b)
```

**Time Series of Daily PRCP**



**Time Series of Daily TMAX**

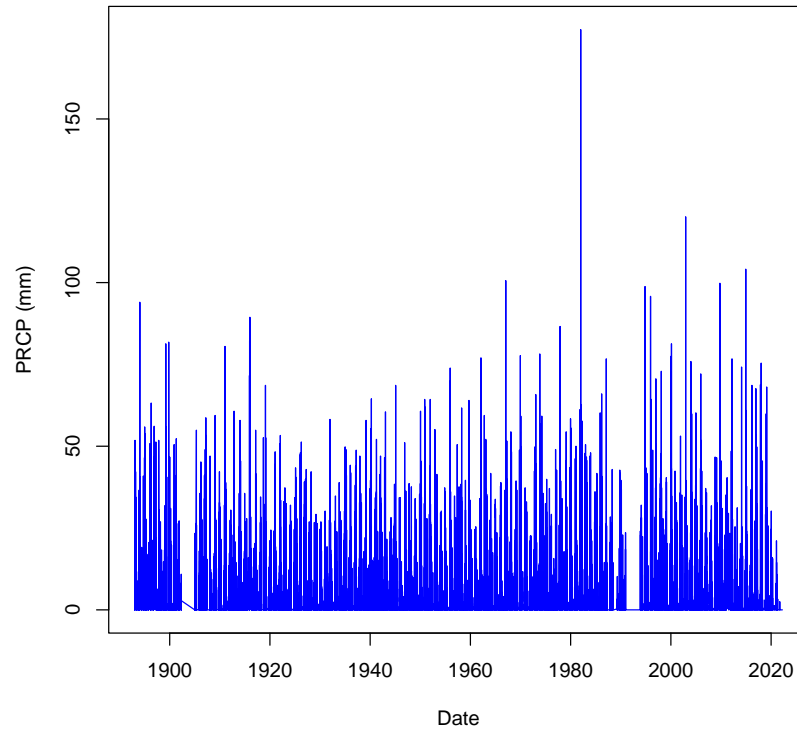


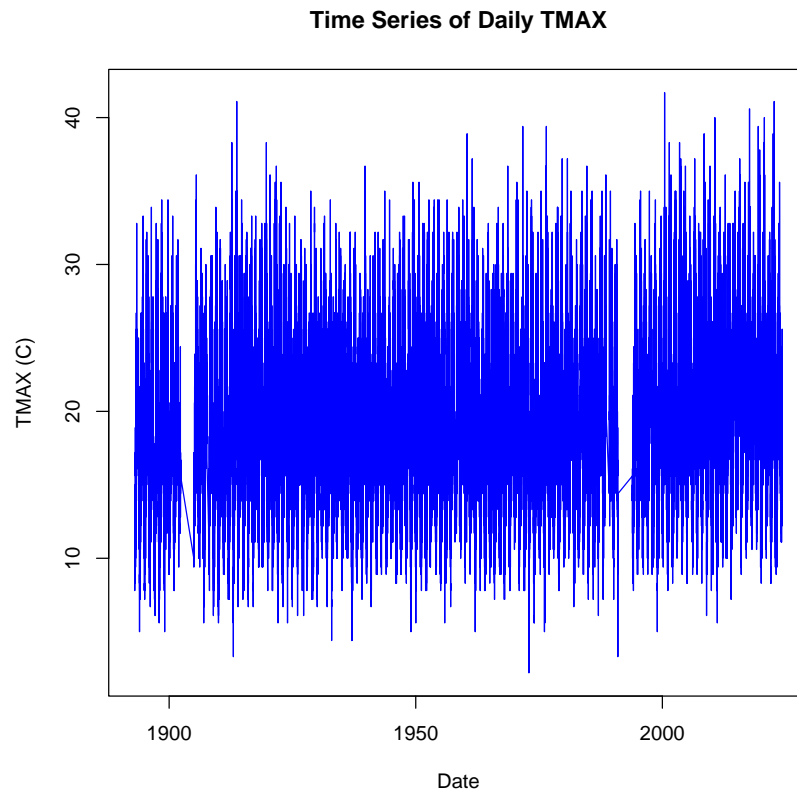


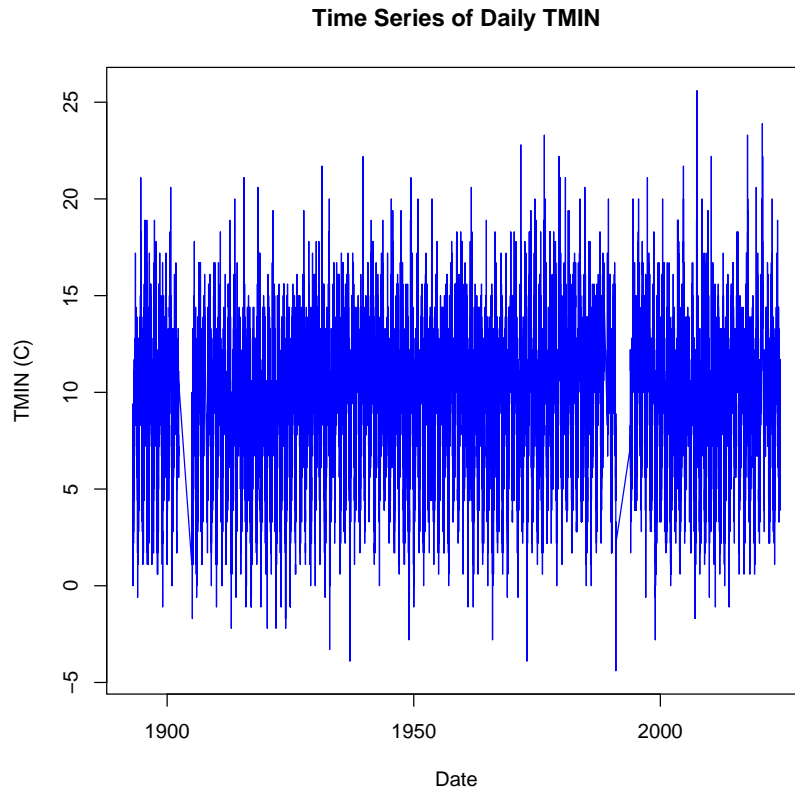
```
QAQC.fun(USC00040693b)
```



**Time Series of Daily PRCP**







If you have any hints of data problems (wacky values, odd patterns), let's sort them out together!

**Function to Create Monthly Values and Normals** For TMAX and TMIN, we want monthly means, for rainfall, we'll want monthly totals.<sup>1</sup>

**Function:** `MonthlyValues.fun()`

**Function:** `NormalValues.fun()`

Example of how to use the functions:

```
USC00042294.monthly <- MonthlyValues.fun(USC00042294b)
USC00042294.normals <- MonthlyNormals.fun(USC00042294b)
USC00040693.monthly <- MonthlyValues.fun(USC00040693b)
USC00040693.normals <- MonthlyNormals.fun(USC00040693b)
```

---

<sup>1</sup>Brody: We need code to exclude months with missing data, these might not be representative of the month if missing, especially for PRCP!

**Function to Create Anomalies** Example of how to use the function:

```
USC00042294.anomalies <-  
  MonthlyAnomalies.fun(USC00042294.monthly, USC00042294.normals)  
USC00040693.anomalies <-  
  MonthlyAnomalies.fun(USC00040693.monthly, USC00040693.normals)
```

## 2.6 Checking on the Results

You can double check that the dataframes you have been making are actually present by using the `ls()` function the console again.

```
ls()  
  
## [1] "coverage.fun"          "datafolder"          "df_exists"  
## [4] "df_head"              "df_names"            "fixDates.fun"  
## [7] "fixValues.fun"         "MonthlyAnomalies.fun" "MonthlyNormals.fun"  
## [10] "MonthlyValues.fun"     "QAQC.fun"            "ReadStations2.fun"  
## [13] "SaveCleanUp.fun"       "sortStations.fun"    "USC00040693"  
## [16] "USC00040693.anomalies" "USC00040693.monthly" "USC00040693.normals"  
## [19] "USC00040693a"         "USC00040693b"        "USC00041614"  
## [22] "USC00042294"          "USC00042294.anomalies" "USC00042294.monthly"  
## [25] "USC00042294.normals"  "USC00042294a"        "USC00042294b"  
## [28] "USC00043157"          "USC00044259"         "USC00044412"  
## [31] "USC00046074"          "USC00046136"         "USC00046719"  
## [34] "USC00046826"          "USC00047821"         "USC00047902"  
## [37] "USC00048351"          "USC00049866"         "USW00023271"
```

Getting into the data is a bit tricky. The datasets is a list of dataframes. Each dataframe is a different variable, where 1 is TMAX, 2 is TMIN, and 3 is PRCP.

The get access to each you use the following code:

- `USC00042294.anomalies[[1]]` for TMAX
- `USC00042294.anomalies[[2]]` for TMIN
- `USC00042294.anomalies[[3]]` for PRCP

Your are done with this guide, now take a break, a walk, and enjoy some screen free downtime. Next, go to Guide 3!

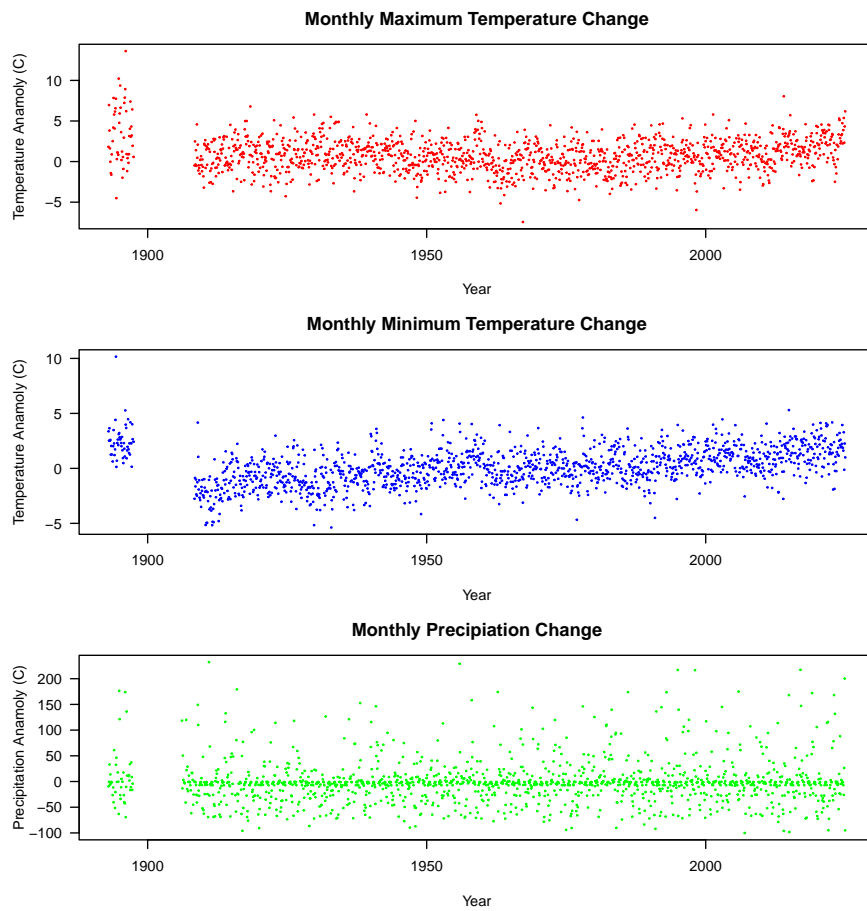


Figure 1: Here's a quick test to see if my data are coming out as expected and as a preview.

## 2.7 Clean Up R Environment

I tend to avoid having lots of objects in my R environment. I like to clean up after myself.

Also, I like to move the final products into csv files. Here's the function to do that. However, I noticed that it's deleting needed files. YIKES! It's also clunky, but not that important at this point, nevertheless, I'll work on it later. So, I suggest you don't run this code until I fix it.

```
CleanUp.fun(datafolder, USC00042294.anamolies, "USC00042294")
```

For now, I am just saving the anamolies data into an RData file that I can use in other Rmd Guides by loading. I doubt you'll need to do that if you use only one Rmd file to knit all the functions together.

```
SaveCleanUp.fun(datafolder)
```

## 3 Describing Marc's Custom Functions

### 3.1 ReadStations2.fun

```
## function (datafolder)
## {
##   my.stations = read.csv(paste0(datafolder, "my.inventory.csv"))
##   colnames <- c("ID", "DATE", "ELEMENT", "VALUE", "M-FLAG",
##               "Q-FLAG", "S-FLAG", "OBS-TIME")
##   for (i in 1:nrow(my.stations)) {
##     assign(my.stations$ID[i], read.csv(paste0(datafolder,
##         noquote(my.stations$ID[i]), ".csv"), header = TRUE,
##         col.names = colnames), envir = parent.frame())
##   }
## }
## <bytecode: 0x2772c58>
```

This is a paragraph. Take 2?

#### 3.1.1 Trouble Shooting

The function fails to read some csv files into the R environment for a variety of reasons that I haven't been able to solve.

To work around the issues, we'll have to read the csv files directly. Here's how an example of how we might do it:

```
USC00042294 <- read.csv("data/USC00042294.csv")
```

### 3.2 fixdates.fun

```
## function (station)
## {
##   station$Ymd = as.Date(as.character(station$DATE), format = "%Y%m%d")
##   station$MONTH = as.numeric(format(station$Ymd, "%m"))
##   station$YEAR = as.numeric(format(station$Ymd, "%Y"))
##   return(station)
## }
## <bytecode: 0x3e83e58>
```

### 3.3 ConvertUnits.fun

```
## function (station)
## {
##     station$VALUE = station$VALUE/10
##     return(station)
## }
## <bytecode: 0x386f3d0>
```

### 3.4 QAQC.fun

```
## function (station)
## {
##     par(mfrow = c(1, 1))
##     plot(VALUE ~ Ymd, data = subset(station, subset = ELEMENT ==
##         "PRCP"), type = "l", col = "blue", main = "Time Series of Daily PRCP",
##         xlab = "Date", ylab = "PRCP (mm)")
##     plot(VALUE ~ Ymd, data = subset(station, subset = ELEMENT ==
##         "TMAX"), type = "l", col = "blue", main = "Time Series of Daily TMAX",
##         xlab = "Date", ylab = "TMAX (C)")
##     plot(VALUE ~ Ymd, data = subset(station, subset = ELEMENT ==
##         "TMIN"), type = "l", col = "blue", main = "Time Series of Daily TMIN",
##         xlab = "Date", ylab = "TMIN (C)")
##     station = subset(station, Q.FLAG != "")
##     station = subset(station, M.FLAG != "")
##     station = subset(station, S.FLAG != "")
##     return(station)
## }
## <bytecode: 0xd381668>
```

#### 3.4.1 How to Evaluate QA/QC Problems

#### 3.4.2 QA/QC Trouble Shooting

I will be getting all the guides working before working on this! But if there are errors with the custom function, this is where workarounds will be described! Please Slack me and mentors if you have any problems!

#### 3.4.3 Coverage Problems (<95%)

```
## function (station, element = "TMAX")
## {
##     Dates.all = data.frame(Ymd = seq.Date(from = min(station$Ymd),
##         to = max(station$Ymd), by = "day"))
```



```
##      station.full = merge(Dates.all, station, all = TRUE)
##      station.coverage = sum(!is.na(station.full$VALUE[station.full$ELEMENT ==
##          element]))/length(station.full$VALUE[station.full$ELEMENT ==
##          element]) * 100
##      return(round(station.coverage, 2))
##  }
## <bytecode: 0x344f600>
```

If you have too many stations that don't have enough data, we'll need to download additional stations. I can show you a way to decipher that ahead of time if you'd like to know. Otherwise, I suggest you double the number of stations selected in the code that generated my.inventory. I have increase the number of stations per state to 15, so perhaps that will give you enough to work with.

If you get an error with this function, be sure you are using the correct file – in our case, it should end with an "a". If you are using the wrong file, you'll get an error.

### 3.4.4 Missing Data

Similar to the problem above, missing data can be a problem. The real issue that I can see is that rainfall is so central because we use the data to calculate monthly totals, but if the month is missing data, then we have a severe bias. Same issue in temperature, but because we are looking at averages, it's less likely to be a major source of bias. But we should should check!

## 3.5 Failure to Read csv files

I noticed that the read

### 3.5.1 Plot Anomaly

Graphic has lots of issues. more next time! But here's a start.

```
options(scipen=5)
par(mar=c(4,6,2,5))

plot(TMAX.a ~ YEAR, data = subset(, MONTH == 1),
     las=1, pch=19, col = "blue", cex=.5, #xlab = "Year",
     ylab = "Maximum Temp Anomaly (C)",
     main="January Maximum Temp Anomaly")
mtext("Maximum Temp Anomaly (C)", side = 2, line = 3)
temp.lm = lm(ANOMALY ~ YEAR, data = subset(station1.TMAX, MONTH == 1))
abline(coef(temp.lm), col = "red")
```

We can see huge periods of time where no data was collected. Yikes! I don't think I can use this station.

My custom functions are probably sensitive to missing values, need to work on that!

### 3.6 MonthlyValues.fun

Here's the function for Monthly Normals:

```
## function (x)
## {
##     x.normals = subset(x, Ymd >= "1961-01-01" & Ymd <= "1990-12-31")
##     x.TMAX.normals.monthly = aggregate(VALUE ~ MONTH, data = subset(x.normals,
##         ELEMENT == "TMAX"), mean)
##     names(x.TMAX.normals.monthly) <- c("MONTH", "NORMALS")
##     x.TMIN.normals.monthly = aggregate(VALUE ~ MONTH, data = subset(x.normals,
##         ELEMENT == "TMIN"), mean)
##     names(x.TMIN.normals.monthly) <- c("MONTH", "NORMALS")
##     x.PRCP.normals.month.year = aggregate(VALUE ~ MONTH + YEAR,
##         data = subset(x.normals, ELEMENT == "PRCP"), sum)
##     x.PRCP.normals.monthly = aggregate(VALUE ~ MONTH, data = subset(x.PRCP.normals.month,
##         mean)
##     names(x.PRCP.normals.monthly) <- c("MONTH", "NORMALS")
##     return(list(x.TMAX.normals.monthly, x.TMIN.normals.monthly,
##         x.PRCP.normals.monthly))
## }
## <bytecode: 0x34dbbf8>
```

If you try to submit the wrong station file, e.g. 'USC00042294', you'll get an error. You need to submit the file that ends with an "b".

Here's the function for Monthly Values:

```
## function (x)
## {
##     x.TMAX.monthly = aggregate(VALUE ~ MONTH + YEAR, data = subset(x,
##         ELEMENT == "TMAX"), mean)
##     names(x.TMAX.monthly) <- c("MONTH", "YEAR", "TMAX")
##     x.TMIN.monthly = aggregate(VALUE ~ MONTH + YEAR, data = subset(x,
##         ELEMENT == "TMIN"), mean)
##     names(x.TMIN.monthly) <- c("MONTH", "YEAR", "TMIN")
##     x.PRCP.monthly = aggregate(VALUE ~ MONTH + YEAR, data = subset(x,
##         ELEMENT == "PRCP"), sum)
##     names(x.PRCP.monthly) <- c("MONTH", "YEAR", "PRCP")
##     return(list(x.TMAX.monthly, x.TMIN.monthly, x.PRCP.monthly))
## }
## <bytecode: 0x496d160>
```

If you try to submit the wrong station file, e.g. ‘USC00042294’, you’ll get an error. You need to submit the file that ends with an “b”.

### 3.7 MonthlyAnomalies.fun

Here’s the function:

```
## function (station.monthly, station.normals)
## {
##   for (i in seq_along(station.monthly)) {
##     TMAX <- merge(station.monthly[[1]], station.normals[[1]],
##       by = "MONTH")
##     TMAX$TMAX.a = TMAX$TMAX - TMAX$NORMALS
##     TMAX$Ymd = as.Date(paste(TMAX$YEAR, TMAX$MONTH, "01",
##       sep = "-"))
##     TMIN <- merge(station.monthly[[2]], station.normals[[2]],
##       by = "MONTH")
##     TMIN$TMIN.a = TMIN$TMIN - TMIN$NORMALS
##     TMIN$Ymd = as.Date(paste(TMIN$YEAR, TMIN$MONTH, "01",
##       sep = "-"))
##     PRCP <- merge(station.monthly[[3]], station.normals[[3]],
##       by = "MONTH")
##     PRCP$PRCP.a = PRCP$PRCP - PRCP$NORMALS
##     PRCP$Ymd = as.Date(paste(PRCP$YEAR, PRCP$MONTH, "01",
##       sep = "-"))
##     return(list(TMAX = TMAX, TMIN = TMIN, PRCP = PRCP))
##   }
## }
## <bytecode: 0x17a1f4b0>
```

If this function fails, it’s likely because the products from the MonthlyNormals.fun or MonthlyValues.fun functions did not work or the products of these functions are not in the R environment and correctly specified in the function call.

## 4 Trouble Shooting and Work Arouds

### 4.1 Checking Results Step by Step

I generally check every step of the way to make sure the function is working. You can look the “Global Environment” in RStudio to see if the objects are there.

In addition, I have written some custom functions to evaluate the dataframes we have (hopefully) created!

- Does these object exist?

```
df_exists("USC00042294b")
```

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

- What are the names within the dataframe?

```
df_names(USC00042294b)
```

```
## [1] "ID"      "DATE"    "ELEMENT" "VALUE"   "M.FLAG"  "Q.FLAG"
## [7] "S.FLAG" "OBS.TIME" "Ymd"     "MONTH"   "YEAR"
```

- What are the first few rows of the dataframe?

```
df_head(USC00042294.anomalies)
```

```
## $TMAX
##   MONTH YEAR      TMAX  NORMALS      TMAX.a      Ymd
## 1      1 1893 13.335484 11.51748  1.8180081 1893-01-01
## 2      1 1975 11.564516 11.51748  0.0470404 1975-01-01
## 3      1 1965 10.645161 11.51748 -0.8723144 1965-01-01
## 4      1 1911 13.051613 11.51748  1.5341372 1911-01-01
## 5      1 1955  9.912903 11.51748 -1.6045725 1955-01-01
## 6      1 1945 11.935484 11.51748  0.4180081 1945-01-01
##
## $TMIN
##   MONTH YEAR      TMIN  NORMALS      TMIN.a      Ymd
## 1      1 1893  5.754839  2.387271  3.3675679 1893-01-01
## 2      1 1975  1.693548  2.387271 -0.6937224 1975-01-01
## 3      1 1965  3.996774  2.387271  1.6095034 1965-01-01
## 4      1 1911  1.338710  2.387271 -1.0485611 1911-01-01
## 5      1 1955  1.838710  2.387271 -0.5485611 1955-01-01
## 6      1 1945  1.416129  2.387271 -0.9711417 1945-01-01
##
## $PRCP
##   MONTH YEAR  PRCP  NORMALS  PRCP.a      Ymd
## 1      1 1893   90.7   100.21  -9.51 1893-01-01
## 2      1 1974   84.0   100.21 -16.21 1974-01-01
## 3      1 1964  101.8   100.21   1.59 1964-01-01
## 4      1 1913   87.1   100.21 -13.11 1913-01-01
## 5      1 2025    5.1   100.21 -95.11 2025-01-01
## 6      1 1944   72.0   100.21 -28.21 1944-01-01
```

## 5 Next Steps

### 5.1 Apply Function to All Stations

So far, I have only run function for 1 station, but I suspect you can figure out how to run it for each one!

This is all we need to do so far. Next week, we'll look at different way to visualize the data!

### 5.2 Save and Clean Up R Environment

I'll save all the station data into csv files, then use them in the next guide to clean, process, and visualize data. I don't think the function is all that useful, so I can show you better ways of doing thin is class.

```
SaveCleanUp.fun(datafolder)
```

```
station1.clean=cleandataframe.fun(station1)
```