

## ENTERING AND CLEANING DATA #3

# CLEANING VERY MESSY DATA

# HURRICANE TRACKING DATA

One version of Atlantic basin hurricane tracks is available here: <http://www.nhc.noaa.gov/data/hurdat/hurdat2-1851-2015-070616.txt>. The data is not in a classic delimited format:

[illegible]

# HURRICANE TRACKING DATA

This data is formatted in the following way:

- Data for many storms are included in one file.
- Data for a storm starts with a shorter line, with values for the storm ID, name, and number of observations for the storm. These values are comma separated.
- Observations for each storm are longer lines. There are multiple observations for each storm, where each observation gives values like the location and maximum winds for the storm at that time.

# HURRICANE TRACKING DATA

Strategy for reading in very messy data:

- ① Read in all lines individually.
- ② Use regular expressions to split each line into the elements you'd like to use to fill columns.
- ③ Write functions, loops, or `apply` calls to process lines and use the contents to fill a data frame.
- ④ Once you have the data in a data frame, do any remaining cleaning to create a data frame that is easy to use to answer research questions.

# HURRICANE TRACKING DATA

Because the data is not nicely formatted, you can't use `read_csv` or similar functions to read it in.

However, the `readLines` function allows you to read a text file in one line at a time. You can then write code and functions to parse the file one line at a time, to turn it into a dataframe you can use.

# HURRICANE TRACKING DATA

The `readLines` function will read in lines from a text file directly, without trying to separate into columns. You can use the `n` argument to specify the number of lines to read it.

For example, to read in three lines from the hurricane tracking data, you can run:

```
tracks_url <- paste0("http://www.nhc.noaa.gov/data/hurdat/",  
                     "hurdat2-1851-2015-070616.txt")  
hurr_tracks <- readLines(tracks_url, n = 3)  
hurr_tracks
```

```
## [1] "AL011851,          UNNAMED,      14,"  
## [2] "18510625, 0000,    , HU, 28.0N,  94.8W,  80, -999, -999, -  
## [3] "18510625, 0600,    , HU, 28.0N,  95.4W,  80, -999, -999, -"
```

# HURRICANE TRACKING DATA

The data has been read in as a vector, rather than a dataframe:

```
class(hurr_tracks)
```

```
## [1] "character"
```

```
length(hurr_tracks)
```

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

```
hurr_tracks[1]
```

```
## [1] "AL011851,          UNNAMED,          14,"
```



# HURRICANE TRACKING DATA

You can use regular expressions to break each line up. For example, you can use `str_split` from the `stringr` package to break the first line of the hurricane track data into its three separate components:

```
library(stringr)
str_split(hurr_tracks[1], pattern = ",")
```

```
## [[1]]
## [1] "AL011851"      "                UNNAMED" "      14"
## [4] ""
```

# HURRICANE TRACKING DATA

You can use this to create a list where each element of the list has the split-up version of a line of the original data. First, read in all of the data:

```
tracks_url <- paste0("http://www.nhc.noaa.gov/data/hurdat/",  
                     "hurdat2-1851-2015-070616.txt")  
hurr_tracks <- readLines(tracks_url)  
length(hurr_tracks)
```

```
## [1] 50919
```

# HURRICANE TRACKING DATA

Next, use `lapply` with `str_split` to split each line of the data at the commas:

```
hurr_tracks <- lapply(hurr_tracks, str_split,  
                      pattern = ",",  
                      simplify = TRUE)
```

```
hurr_tracks[[1]]
```

```
##      [,1]      [,2]      [,3]      [,4]  
## [1,] "AL011851" "      UNNAMED" "      14" ""
```

```
hurr_tracks[[2]][1:2]
```

```
## [1] "18510625" " 0000"
```

# HURRICANE TRACKING DATA

Next, you want to split this list into two lists, one with the shorter “meta-data” lines and one with the longer “observation” lines. You can use `sapply` to create a vector with the length of each line. You will later use this to identify which lines are short or long.

```
hurr_lengths <- sapply(hurr_tracks, length)
hurr_lengths[1:17]
```

```
## [1] 4 21 21 21 21 21 21 21 21 21 21 21 21 21 21 21 4 21
```

```
unique(hurr_lengths)
```

```
## [1] 4 21
```

# HURRICANE TRACKING DATA

You can use bracket indexing to split the `hurr_tracks` into two lists: one with the shorter lines that start each observation (`hurr_meta`) and one with the storm observations (`hurr_obs`). Use bracket indexing with the `hurr_lengths` vector you just created to make that split.

```
hurr_meta <- hurr_tracks[hurr_lengths == 4]  
hurr_obs <- hurr_tracks[hurr_lengths == 21]
```

# HURRICANE TRACKING DATA

```
hurr_meta[1:3]
```

```
## [[1]]  
##      [,1]      [,2]      [,3]      [,4]  
## [1,] "AL011851" "      " UNNAMED" "      14" ""  
##  
## [[2]]  
##      [,1]      [,2]      [,3]      [,4]  
## [1,] "AL021851" "      " UNNAMED" "      1" ""  
##  
## [[3]]  
##      [,1]      [,2]      [,3]      [,4]  
## [1,] "AL031851" "      " UNNAMED" "      1" ""
```

# HURRICANE TRACKING DATA

```
hurr_obs[1:2]
```

```
## [[1]]
##      [,1]      [,2]      [,3] [,4]  [,5]      [,6]      [,7]
## [1,] "18510625" " 0000" "  " " HU" " 28.0N" " 94.8W" " 80"
##      [,9]      [,10]     [,11]    [,12]    [,13]    [,14]    [,15]
## [1,] " -999" " -999" " -999" " -999" " -999" " -999" " -999"
##      [,17]     [,18]     [,19]    [,20]    [,21]
## [1,] " -999" " -999" " -999" " -999" ""
##
## [[2]]
##      [,1]      [,2]      [,3] [,4]  [,5]      [,6]      [,7]
## [1,] "18510625" " 0600" "  " " HU" " 28.0N" " 95.4W" " 80"
##      [,9]      [,10]     [,11]    [,12]    [,13]    [,14]    [,15]
## [1,] " -999" " -999" " -999" " -999" " -999" " -999" " -999"
##      [,17]     [,18]     [,19]    [,20]    [,21]
## [1,] " -999" " -999" " -999" " -999" ""
```

# HURRICANE TRACKING DATA

Now, you can use `bind_rows` from `dplyr` to change the list of metadata into a dataframe. (You first need to use `as_tibble` with `lapply` to convert all elements of the list from matrices to dataframes.)

```
library(dplyr)
hurr_meta <- lapply(hurr_meta, tibble::as_tibble)
hurr_meta <- bind_rows(hurr_meta)
hurr_meta %>%
  slice(1:3)
```

```
## # A tibble: 3 × 4
##       V1          V2      V3    V4
##   <chr>    <chr>  <chr> <chr>
## 1 AL011851 UNNAMED    14
## 2 AL021851 UNNAMED     1
## 3 AL031851 UNNAMED     1
```



# HURRICANE TRACKING DATA

You can clean up the data a bit more.

- First, the fourth column doesn't have any non-missing values, so you can get rid of it:

```
unique(hurr_meta$V4)
```

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

- Second, the second and third columns include a lot of leading whitespace:

```
hurr_meta$V2[1:2]
```

```
## [1] "                UNNAMED" "                UNNAMED"
```

- Last, we want to name the columns.

# HURRICANE TRACKING DATA

```
hurr_meta <- hurr_meta %>%  
  select(-V4) %>%  
  rename(storm_id = V1, storm_name = V2, n_obs = V3) %>%  
  mutate(storm_name = str_trim(storm_name),  
         n_obs = as.numeric(n_obs))  
hurr_meta %>% slice(1:3)
```

```
## # A tibble: 3 × 3  
##   storm_id storm_name n_obs  
##   <chr>      <chr> <dbl>  
## 1 AL011851  UNNAMED    14  
## 2 AL021851  UNNAMED     1  
## 3 AL031851  UNNAMED     1
```

# HURRICANE TRACKING DATA

Now you can do the same idea with the hurricane observations. First, we'll want to add storm identifiers to that data. The "meta" data includes storm ids and the number of observations per storm. We can take advantage of that to make a `storm_id` vector that will line up with the storm observations.

```
storm_id <- rep(hurr_meta$storm_id, times = hurr_meta$n_obs)
head(storm_id, 3)
```

```
## [1] "AL011851" "AL011851" "AL011851"
```

```
length(storm_id)
```

```
## [1] 49105
```

```
length(hurr_obs)
```

```
## [1] 49105
```

# HURRICANE TRACKING DATA

```
hurr_obs <- lapply(hurr_obs, tibble::as_tibble)
hurr_obs <- bind_rows(hurr_obs) %>%
  mutate(storm_id = storm_id)
hurr_obs %>% select(V1, V2, V5, V6, storm_id) %>% slice(1:3)
```

```
## # A tibble: 3 × 5
```

	V1	V2	V5	V6	storm_id
	<chr>	<chr>	<chr>	<chr>	<chr>
## 1	18510625	0000	28.0N	94.8W	AL011851
## 2	18510625	0600	28.0N	95.4W	AL011851
## 3	18510625	1200	28.0N	96.0W	AL011851

# HURRICANE TRACKING DATA

To finish, you just need to clean up the data. Now that the data is in a dataframe, this process is inline with what you've been doing with `dplyr` and related packages.

The “README” file for the hurricane tracking data is useful at this point:

`http:`

`//www.nhc.noaa.gov/data/hurdat/hurdat2-format-atlantic.pdf`

# HURRICANE TRACKING DATA

First, say you only want some of the columns for a study you are doing. You can use `select` to clean up the dataframe by limiting it to columns you need.

If you only need date, time, storm status, location (latitude and longitude), maximum sustained winds, and minimum pressure, then you can run:

```
hurr_obs <- hurr_obs %>%  
  select(V1, V2, V4:V8, storm_id) %>%  
  rename(date = V1, time = V2, status = V4, latitude = V5,  
         longitude = V6, wind = V7, pressure = V8)  
hurr_obs %>% slice(1:3) %>%  
  select(date, time, status, latitude, longitude)
```

```
## # A tibble: 3 × 5  
##       date   time status latitude longitude  
##   <chr> <chr>   <chr>    <chr>    <chr>  
## 1 18510625 0000    HU    28.0N    94.8W  
## 2 18510625 0600    HU    28.0N    95.4W  
## 3 18510625 1200    HU    28.0N    96.0W
```

# HURRICANE TRACKING DATA

Next, the first two columns give the date and time. You can unite these and then convert them to a Date-time class.

```
library(tidyr)
library(lubridate)
hurr_obs <- hurr_obs %>%
  unite(date_time, date, time) %>%
  mutate(date_time = ymd_hm(date_time))
hurr_obs %>% slice(1:3) %>%
  select(date_time, status, latitude, longitude)
```

```
## # A tibble: 3 × 4
```

##		date_time	status	latitude	longitude
##		<dtm>	<chr>	<chr>	<chr>
## 1	1851-06-25 00:00:00	HU	28.0N	94.8W	
## 2	1851-06-25 06:00:00	HU	28.0N	95.4W	
## 3	1851-06-25 12:00:00	HU	28.0N	96.0W	

# HURRICANE TRACKING DATA

Next, you can change status to a factor and give the levels more meaningful names:

```
unique(hurr_obs$status)
```

```
## [1] " HU" " TS" " EX" " TD" " LO" " DB" " SD" " SS" " WV"
```

```
storm_levels <- c("TD", "TS", "HU", "EX",  
                  "SD", "SS", "LO", "WV", "DB")  
storm_labels <- c("Tropical depression", "Tropical storm",  
                  "Hurricane", "Extratropical cyclone",  
                  "Subtropical depression",  
                  "Subtropical storm", "Other low",  
                  "Tropical wave", "Disturbance")  
hurr_obs <- hurr_obs %>%  
  mutate(status = factor(str_trim(status),  
                          levels = storm_levels,  
                          labels = storm_labels))
```



# HURRICANE TRACKING DATA

Now, you can clean up the latitude and longitude. Ultimately, we'll want numeric values for those so we can use them for mapping. You can use regular expressions to separate the numeric and non-numeric parts of these columns. For example:

```
head(str_extract(hurr_obs$latitude, "[A-Z]"))
```

```
## [1] "N" "N" "N" "N" "N" "N"
```

```
head(str_extract(hurr_obs$latitude, "[^A-Z]+"))
```

```
## [1] " 28.0" " 28.0" " 28.0" " 28.1" " 28.2" " 28.2"
```



# HURRICANE TRACKING DATA

Now these elements are in separate columns:

```
hurr_obs %>%  
  select(latitude, lat_dir, longitude, lon_dir) %>%  
  slice(1:2)
```

```
## # A tibble: 2 × 4  
##   latitude lat_dir longitude lon_dir  
##   <dbl>   <chr>     <dbl>   <chr>  
## 1      28      N       94.8     W  
## 2      28      N       95.4     W
```

```
unique(hurr_obs$lat_dir)
```

```
## [1] "N"
```

```
unique(hurr_obs$lon_dir)
```

```
## [1] "W" "E"
```

# HURRICANE TRACKING DATA

If we're looking at US impacts, we probably only need observations from the western hemisphere, so let's filter out other values:

```
hurr_obs <- hurr_obs %>%  
  filter(lon_dir == "W")
```

# HURRICANE TRACKING DATA

Next, clean up the wind column:

```
unique(hurr_obs$wind)[1:5]
```

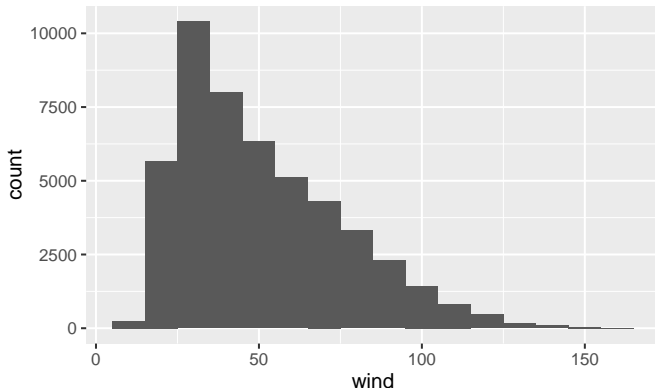
```
## [1] " 80" " 70" " 60" " 50" " 40"
```

```
hurr_obs <- hurr_obs %>%  
  mutate(wind = ifelse(wind == "-99", NA,  
                        as.numeric(wind)))
```

# HURRICANE TRACKING DATA

Check the cleaned measurements:

```
library(ggplot2)
ggplot(hurr_obs, aes(x = wind)) +
  geom_histogram(binwidth = 10)
```



# HURRICANE TRACKING DATA

Clean and check air pressure measurements in the same way:

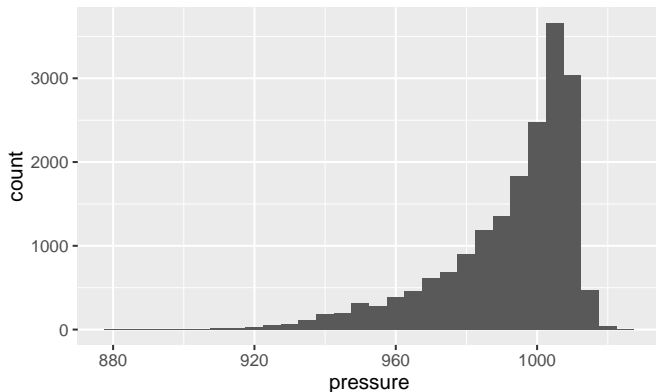
```
head(unique(hurr_obs$pressure))
```

```
## [1] " -999" " 961" " 924" " 938" " 950" " 997"
```

```
hurr_obs <- hurr_obs %>%  
  mutate(pressure = ifelse(pressure == " -999", NA,  
                           as.numeric(pressure)))
```

# HURRICANE TRACKING DATA

```
ggplot(hurr_obs, aes(x = pressure)) +  
  geom_histogram(binwidth = 5)
```





# HURRICANE TRACKING DATA

Check some of the very low pressure measurements:

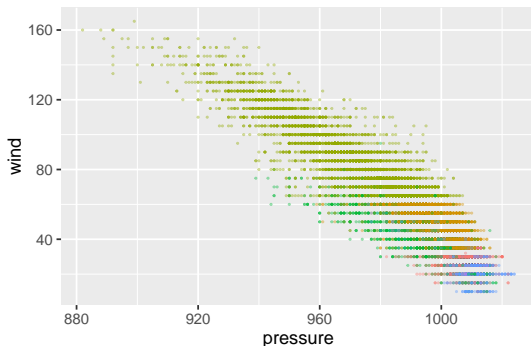
```
hurr_obs %>% arrange(pressure) %>%  
  select(date_time, wind, pressure) %>% slice(1:5)
```

```
## # A tibble: 5 × 3  
##           date_time  wind pressure  
##           <dtm>    <dbl>    <dbl>  
## 1 2005-10-19 12:00:00   160     882  
## 2 1988-09-14 00:00:00   160     888  
## 3 1988-09-14 06:00:00   155     889  
## 4 1935-09-03 00:00:00   160     892  
## 5 1935-09-03 02:00:00   160     892
```

# HURRICANE TRACKING DATA

Explore pressure versus wind speed, by storm status:

```
ggplot(hurr_obs, aes(x = pressure, y = wind,  
                     color = status)) +  
  geom_point(size = 0.2, alpha = 0.4)
```



status

- Tropical depression
- Tropical storm
- Hurricane
- Extratropical cyclone
- Subtropical depression
- Subtropical storm
- Other low
- Tropical wave
- Disturbance

# HURRICANE TRACKING DATA

Next, we want to map storms by decade. Add hurricane decade:

```
hurr_obs <- hurr_obs %>%  
  mutate(decade = substring(year(date_time), 1, 3),  
         decade = paste0(decade, "0s"))  
unique(hurr_obs$decade)
```

```
## [1] "1850s" "1860s" "1870s" "1880s" "1890s" "1900s" "1910s"  
## [9] "1930s" "1940s" "1950s" "1960s" "1970s" "1980s" "1990s"  
## [17] "2010s"
```

Add logical for whether the storm was ever category 5:

```
hurr_obs <- hurr_obs %>%  
  group_by(storm_id) %>%  
  mutate(cat_5 = max(wind) >= 137) %>%  
  ungroup()
```

# HURRICANE TRACKING DATA

To map the hurricane tracks, you need a base map to add the tracks to.  
Pull data to map hurricane-prone states:

```
east_states <- c("florida", "georgia", "south carolina",  
                "north carolina", "virginia", "maryland",  
                "delaware", "new jersey", "new york",  
                "connecticut", "massachusetts",  
                "rhode island", "vermont", "new hampshire",  
                "maine", "pennsylvania", "west virginia",  
                "tennessee", "kentucky", "alabama",  
                "arkansas", "texas", "mississippi",  
                "louisiana")  
east_us <- map_data("state", region = east_states)
```

# HURRICANE TRACKING DATA

Plot tracks over a map of hurricane-prone states. Add thicker lines for storms that were category 5 at least once in their history.

```
ggplot(east_us, aes(x = long, y = lat, group = group)) +  
  geom_polygon(fill = "cornsilk", color = "cornsilk") +  
  theme_void() +  
  xlim(c(-108, -65)) + ylim(c(23, 48)) +  
  geom_path(data = hurr_obs,  
            aes(x = -longitude, y = latitude,  
                group = storm_id),  
            color = "red", alpha = 0.2, size = 0.2) +  
  geom_path(data = filter(hurr_obs, cat_5),  
            aes(x = -longitude, y = latitude,  
                group = storm_id),  
            color = "red") +  
  facet_wrap(~ decade)
```

# HURRICANE TRACKING DATA

Check trends in maximum wind recorded in any observation each year:

1850s



1860s



1870s



1880s



1890s



1900s



1910s



1920s



1930s



1940s



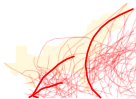
1950s



1960s



1970s



1980s



1990s



2000s



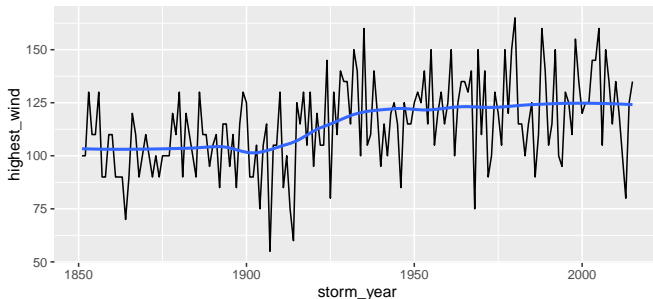
2010s



# HURRICANE TRACKING DATA

Maximum wind observed each year:

```
hurr_obs %>%  
  mutate(storm_year = year(date_time)) %>%  
  group_by(storm_year) %>%  
  summarize(highest_wind = max(wind, na.rm = TRUE)) %>%  
  ggplot(aes(x = storm_year, y = highest_wind)) +  
  geom_line() + geom_smooth(se = FALSE, span = 0.5)
```



# HURRICANE TRACKING DATA

There is an R package named `gender` that predicts whether a name is male or female based on historical data:

Vignette for `gender` package

This package uses one of several databases of names (here, we'll use Social Security Administration data), inputs a year or range of years, and outputs whether a name in that year was more likely female or male.

We can apply a function from this package across all the named storms to see how male / female proportions changed over time.



# HURRICANE TRACKING DATA

First, install the package (as well as `genderdata`, which is required to use the package). Once you do, you can use `gender` to determine the most common gender associated with a name in a given year or range of years:

```
# install.packages("gender")
# install.packages("genderdata")
library(gender)
gender("KATRINA", years = 2005)[ , c("name", "gender")]
```

```
## # A tibble: 1 × 2
##   name gender
##   <chr>  <chr>
## 1 KATRINA female
```

# HURRICANE TRACKING DATA

To apply this function across all our storms, it helps if we write a small function that “wraps” the gender function and outputs exactly (and only) what we want, in the format we want:

```
get_gender <- function(storm_name, storm_year){  
  storm_gender <- gender(names = storm_name,  
                           years = storm_year,  
                           method = "ssa")$gender  
  if(length(storm_gender) == 0) storm_gender <- NA  
  return(storm_gender)  
}
```

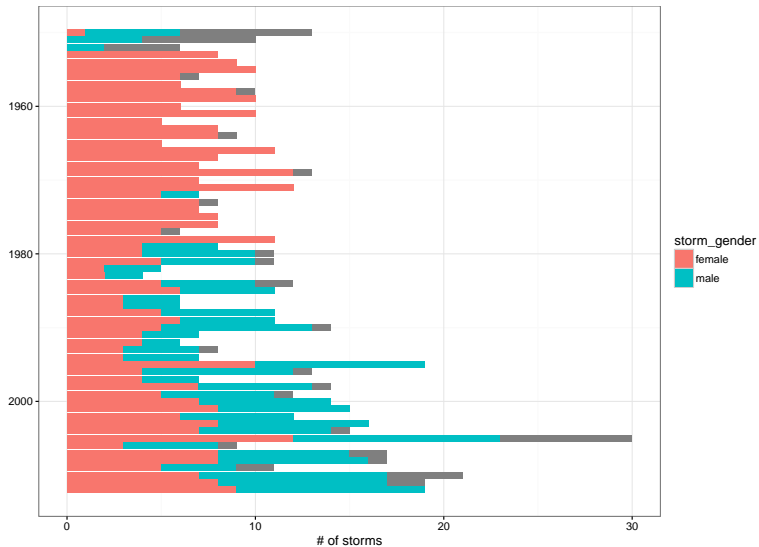


# HURRICANE TRACKING DATA

Now, plot a bar chart with the number of male, female, and unclear storms each year:

```
hurr_genders %>%  
  group_by(storm_year, storm_gender) %>%  
  summarize(n = n()) %>%  
  ggplot(aes(x = storm_year, y = n, fill = storm_gender)) +  
  geom_bar(stat = "identity") +  
  coord_flip() +  
  scale_x_reverse() +  
  theme_bw() +  
  xlab("") + ylab("# of storms")
```

# HURRICANE TRACKING DATA



# HURRICANE TRACKING DATA

Next, you can write a function to plot the track for a specific storm. You'll want to be able to call the function by storm name and year, so join in the storm names from the `hurr_meta` dataset. We'll exclude any "UNNAMED" storms.

```
hurr_obs <- hurr_obs %>%  
  left_join(hurr_meta, by = "storm_id") %>%  
  filter(storm_name != "UNNAMED") %>%  
  mutate(storm_year = year(date_time))
```

Next, write a function to plot the track for a single storm. Use color to show storm status and size to show wind speed.

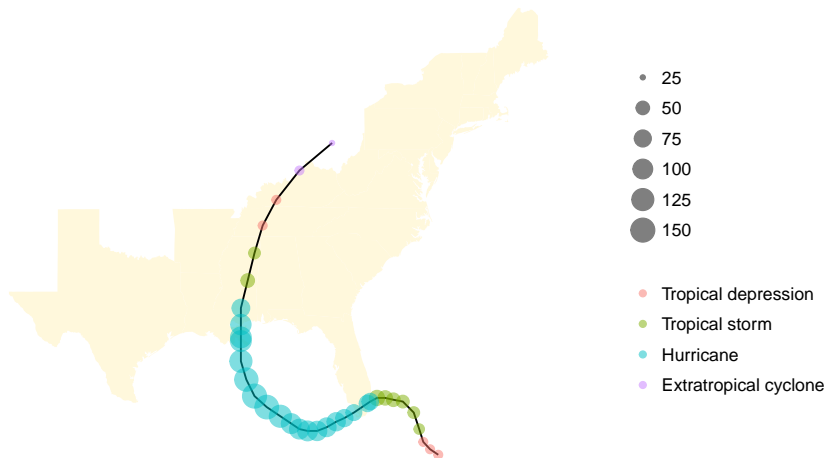
# HURRICANE TRACKING DATA

```
map_track <- function(storm, year, map_data = east_us,
                      hurr_data = hurr_obs){
  to_plot <- hurr_obs %>%
    filter(storm_name == toupper(storm) & storm_year == year)
  out <- ggplot(east_us, aes(x = long, y = lat,
                           group = group)) +
    geom_polygon(fill = "cornsilk") +
    theme_void() +
    xlim(c(-108, -65)) + ylim(c(23, 48)) +
    geom_path(data = to_plot,
              aes(x = -longitude, y = latitude,
                  group = NULL)) +
    geom_point(data = to_plot,
               aes(x = -longitude, y = latitude,
                   group = NULL, color = status,
                   size = wind), alpha = 0.5)

  return(out)
}
```

# HURRICANE TRACKING DATA

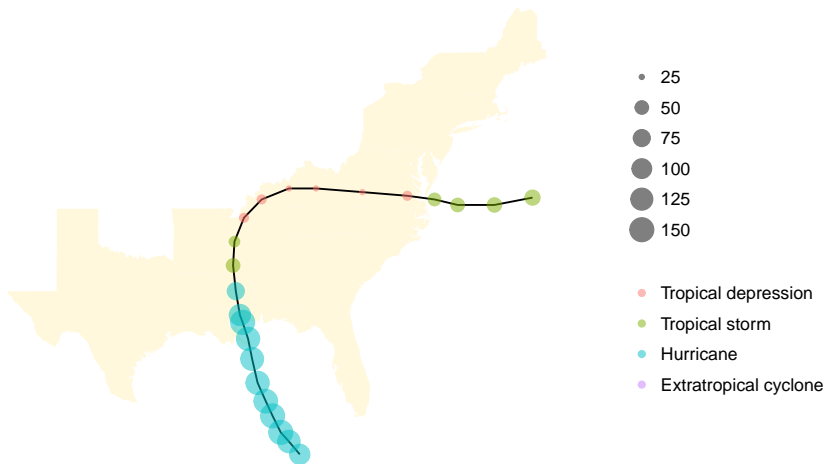
```
map_track(storm = "Katrina", year = "2005")
```





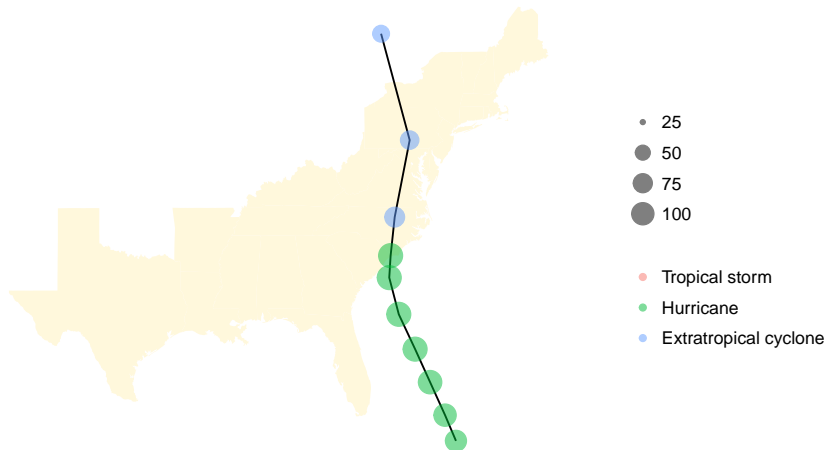
# HURRICANE TRACKING DATA

```
map_track(storm = "Camille", year = "1969")
```



# HURRICANE TRACKING DATA

```
map_track(storm = "Hazel", year = "1954")
```



## READLINES

You can also write code with `readLines` that will read, check, and clean each line, one line at a time.

```
con <- file("~/my_file.txt", open = "r")
while (length(single_line <-
                readLines(con, n = 1,
                          warn = FALSE)) > 0) {

  ## Code to check and clean each line and
  ## then add it to "cleaned" data frame.
  ## Run operations on `single_line`.

}
close(con)
```

This can be particularly useful if you're cleaning a very big file, especially if there are many lines you don't want to keep.

## PULLING ONLINE DATA

# APIs

# EXAMPLE R PACKAGES FOR APIs

# ROpenSci

# PARSING WEBPAGES