

# Cleaning Data

## Cleaning Data In R

GitHub Repository with this demonstration: <https://github.com/unmrds/R-data-cleaning>

Zipfile Download: <https://github.com/unmrds/R-data-cleaning/archive/master.zip>

Check out the **Preflight Check** to see if you have the needed R libraries installed run the `package.check.R` script in the top directory of this R project.

When planning a data analysis the first step, and often most time consuming, is the acquisition and processing of the data into a form that can be used in the analytic procedures you intend to use. Today we are going to focus on a sequence of steps that generally follow the workflow that you will find yourself going through when bringing data into R to perform an analysis.

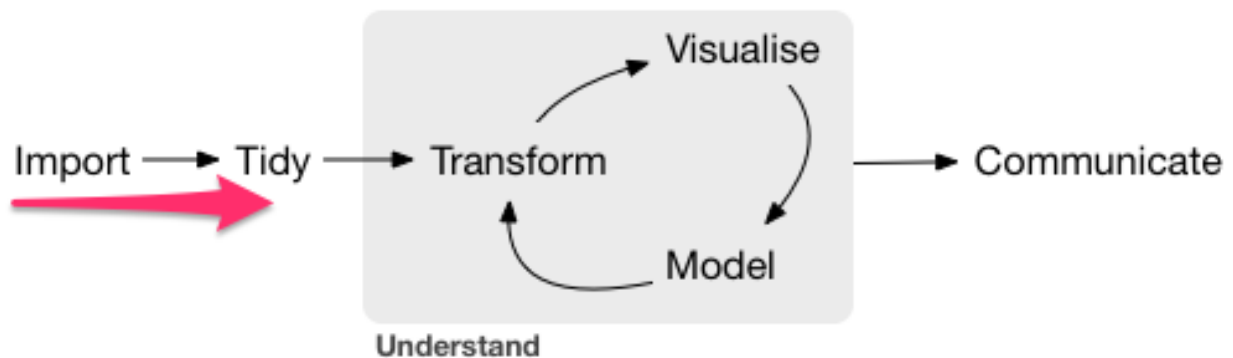


Figure 1: Portion of *R for Data Science*<sup>1</sup> workflow

### Deal with issues that may come up when importing data files

1. Identify and correct structural issues in the source data that prevent clean import into R data structures
2. Check and handle data type errors
3. Check and handle missing data values

### Tuning up the structure of the data to facilitate analysis

4. Split up fields that contain multiple values in a single field
5. Check for anomalous values and otherwise explore the data to become familiar with its content and structure.

*Beyond what we will cover today* - continued structural changes and the rest of the exploration, analysis, and communication process.

### Data for today's demonstration

The data for this demonstration are based upon the `idigbio_rodents.csv` dataset. The data are described as follows in the repository where they are shared:

<sup>1</sup>Hadley Wickham & Garrett Golemund. 2017. *R for Data Science*. O'Reilly. <https://r4ds.had.co.nz>

The `idigbio_rodents.csv` dataset is just over 10k records and contains data from natural history collections specimen records representing 5 genera from 4 major US collections, limited to US records. All records are from the Order Rodentia. All the data are mapped to the biodiversity data standard Darwin Core (<http://rs.tdwg.org/dwc/terms/>).

The original data have been modified for use in this demonstration by:

1. Generating new data columns (`latDMS` and `lonDMS`) for latitude and longitude that have sample coordinates presented in Degrees-Minutes-Seconds instead of the originally provided decimal degrees.
2. Generating a column of mixed numeric and text values - `textLatDD`.

This is the `../data/learning.csv` file. These newly created columns in addition to some of the originally provided ones will be used to demonstrate a variety of data cleaning steps in R.

An additional file was developed that only includes the first 10 rows of the file (including headers) but introduces a structural error. This file is the `../data/learning_struct.csv` file.

## R libraries used in the demonstration

For this demonstration a set of R packages that are part of the *Tidyverse*. The tidyverse collection of packages provide (as described on the project's homepage):

an opinionated collection of R packages designed for data science. All packages share an underlying design philosophy, grammar, and data structures.

There are currently over 14,000 R packages in the Comprehensive R Archive Network (CRAN). While the tidyverse packages provide a useful degree of consistency and internal interoperability it is strongly encouraged to examine the broad collection of R packages when working on a particular analysis problem.

If you need to install the tidyverse packages in your environment you can execute the `install.packages("tidyverse")` command.

## 1. Identify and correct structural issues in the source data that prevent clean import into R data structures

R can import a wide variety of *rectangular* data structures: comma-delimited, tab-delimited, excel spreadsheets, fixed-width among the many options. If there are errors in the structure of these files, R import commands may not be able to parse the lines in the data file preventing import. In these cases the returned error messages may provide some clues to where the errors may be found.

**One strategy for identifying potential structural issues in the source file is to try to import the dataset and review any errors that are returned**

Let's try it first with a small file ...

```
library(tidyverse)

## -- Attaching packages -----
## v ggplot2 3.2.1    v purrr  0.3.3
## v tibble  2.1.3    v dplyr  0.8.4
## v tidyr   1.0.2    v stringr 1.4.0
## v readr   1.3.1    v forcats 0.4.0

## -- Conflicts -----
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()

# Import the source CSV file that contains a structural flaw
rawDataStruct <- read_csv("../data/learning_struct.csv",
                           progress = FALSE)
```

```

## Warning: Missing column names filled in: 'X1' [1]

## Parsed with column specification:
## cols(
##   .default = col_character(),
##   X1 = col_double(),
##   catalogNumber = col_double(),
##   year = col_double(),
##   day = col_double()
## )

## See spec(...) for full column specifications.

## Warning: 1 parsing failure.
## row col   expected      actual      file
##   7  -- 24 columns 26 columns '../data/learning_struct.csv'

# Display the column definitions for the imported dataset
spec(rawDataStruct)

## cols(
##   X1 = col_double(),
##   uuid = col_character(),
##   institutionCode = col_character(),
##   collectionCode = col_character(),
##   catalogNumber = col_double(),
##   recordedBy = col_character(),
##   countryCode = col_character(),
##   stateProvince = col_character(),
##   county = col_character(),
##   decimalLatitude = col_character(),
##   decimalLongitude = col_character(),
##   eventDate = col_character(),
##   year = col_double(),
##   month = col_character(),
##   day = col_double(),
##   genus = col_character(),
##   specificEpithet = col_character(),
##   scientificName = col_character(),
##   weight = col_character(),
##   length = col_character(),
##   sex = col_character(),
##   latDMS = col_character(),
##   lonDMS = col_character(),
##   textLatDD = col_character()
## )

# Report the problems that were encountered when the data were imported.
problems(rawDataStruct)

## # A tibble: 1 x 5
##   row col   expected      actual      file
##   <int> <chr> <chr>      <chr>      <chr>
## 1     7 <NA>  24 columns 26 columns '../data/learning_struct.csv'

# Display the imported table
rawDataStruct

```

```
## # A tibble: 9 x 24
##       X1 uuid  institutionCode collectionCode catalogNumber recordedBy
##   <dbl> <chr> <chr>           <chr>           <dbl> <chr>
## 1     1 0603~ mvz             mammal specim~      219088 collector~
## 2     2 0fb1~ mvz             mammal specim~      233524 collector~
## 3     3 1a69~ mvz             mammal specim~      234346 collector~
## 4     4 1a99~ mvz             mammal specim~      233951 collector~
## 5     5 1f3b~ mvz             mammal specim~      235290 collector~
## 6     6 203f~ uam             mammal specim~        85106 collector~
## 7     7 21ce~ omnh            mammals              50048 caldwell
## 8     8 23d3~ mvz             mammal specim~      216309 collector~
## 9     9 244b~ msb             mammal specim~      294933 collector~
## # ... with 18 more variables: countryCode <chr>, stateProvince <chr>,
## #   county <chr>, decimalLatitude <chr>, decimalLongitude <chr>,
## #   eventDate <chr>, year <dbl>, month <chr>, day <dbl>, genus <chr>,
## #   specificEpithet <chr>, scientificName <chr>, weight <chr>, length <chr>,
## #   sex <chr>, latDMS <chr>, lonDMS <chr>, textLatDD <chr>
```

Let's take a look at the source data file and see if we can find the problem...

Now let's try it with the full dataset for which the structural problem has been resolved.

```
library(tidyverse)

# Import the source CSV file that does not contain the structural problem highlighted above
rawData <- read_csv("../data/learning.csv",
                    progress = FALSE)

## Warning: Missing column names filled in: 'X1' [1]

## Parsed with column specification:
## cols(
##   .default = col_character(),
##   X1 = col_double(),
##   catalogNumber = col_double(),
##   decimalLatitude = col_double(),
##   decimalLongitude = col_double(),
##   eventDate = col_datetime(format = ""),
##   year = col_double(),
##   month = col_double(),
##   day = col_double(),
##   weight = col_double(),
##   length = col_double()
## )

## See spec(...) for full column specifications.

# Display the column definitions for the imported dataset
spec(rawData)

## cols(
##   X1 = col_double(),
##   uuid = col_character(),
##   institutionCode = col_character(),
##   collectionCode = col_character(),
##   catalogNumber = col_double(),
##   recordedBy = col_character(),
```

```
## countryCode = col_character(),
## stateProvince = col_character(),
## county = col_character(),
## decimalLatitude = col_double(),
## decimalLongitude = col_double(),
## eventDate = col_datetime(format = ""),
## year = col_double(),
## month = col_double(),
## day = col_double(),
## genus = col_character(),
## specificEpithet = col_character(),
## scientificName = col_character(),
## weight = col_double(),
## length = col_double(),
## sex = col_character(),
## latDMS = col_character(),
## lonDMS = col_character(),
## textLatDD = col_character()
## )

# Report the problems that were encountered when the data were imported.
problems(rawData)
```

```
## [1] row      col      expected actual
## <0 rows> (or 0-length row.names)
```

```
# Display the imported table
rawData
```

```
## # A tibble: 10,767 x 24
##       X1 uuid  institutionCode collectionCode catalogNumber recordedBy
##   <dbl> <chr> <chr>           <chr>           <dbl> <chr>
## 1     1 0603~ mvz             mammal specim~    219088 collector~
## 2     2 0fb1~ mvz             mammal specim~    233524 collector~
## 3     3 1a69~ mvz             mammal specim~    234346 collector~
## 4     4 1a99~ mvz             mammal specim~    233951 collector~
## 5     5 1f3b~ mvz             mammal specim~    235290 collector~
## 6     6 203f~ uam             mammal specim~     85106 collector~
## 7     7 21ce~ omnh            mammals              50048 caldwell,~
## 8     8 23d3~ mvz             mammal specim~    216309 collector~
## 9     9 244b~ msb             mammal specim~    294933 collector~
## 10    10 2682~ uam             mammal specim~     50255 collector~
## # ... with 10,757 more rows, and 18 more variables: countryCode <chr>,
## #   stateProvince <chr>, county <chr>, decimalLatitude <dbl>,
## #   decimalLongitude <dbl>, eventDate <dtm>, year <dbl>, month <dbl>,
## #   day <dbl>, genus <chr>, specificEpithet <chr>, scientificName <chr>,
## #   weight <dbl>, length <dbl>, sex <chr>, latDMS <chr>, lonDMS <chr>,
## #   textLatDD <chr>
```

Some questions:

1. How do the data types for the columns from this import process differ from those in the previous subset (at least before we fixed it)? Why do you think this is the case?
2. Where there any errors identified during the import? If no, does this mean that there are no potential problems or issues with the imported data? **Let's take a look**
3. How would you explain the values in the `eventDate` column when compared to the `year`, `month`, and `day` columns?

4. What were the different ways in which missing data values were handled?

## 2. Checking and handling data type errors

Depending on the types of data that are encountered in each column of the imported dataset different R import functions will automatically “type” the column (or in R terminology set the “mode” of the column) based on some sample of rows from the source file. In the case of `readr`, the data reading package used by tidyverse, the first 1000 lines of data will be read to determine the data type that should be used for each column. The core R data types are:

- character
- numeric (real or decimal)
- integer
- logical
- complex

These core data types can then be used as the foundation for more complex data types such as *dates*, *times* and *datetimes*.

The core data structures that can be used to organize collections of these data types include:

- vector - a sequence of data values of the same type
- list - a sequence of data values of the same or different types, and structures
- matrix - a vector for which *dimensions*
- data frame - a structured collection of vectors of the same length
- factors - a vector where each value is a set of integer values that are associated with a collection of categorical character values

Let's focus on the `catalogNumber` and `textLatDD` columns in our sample datasets.

```
spec(rawData)
```

```
## cols(  
##   X1 = col_double(),  
##   uuid = col_character(),  
##   institutionCode = col_character(),  
##   collectionCode = col_character(),  
##   catalogNumber = col_double(),  
##   recordedBy = col_character(),  
##   countryCode = col_character(),  
##   stateProvince = col_character(),  
##   county = col_character(),  
##   decimalLatitude = col_double(),  
##   decimalLongitude = col_double(),  
##   eventDate = col_datetime(format = ""),  
##   year = col_double(),  
##   month = col_double(),  
##   day = col_double(),  
##   genus = col_character(),  
##   specificEpithet = col_character(),  
##   scientificName = col_character(),  
##   weight = col_double(),  
##   length = col_double(),  
##   sex = col_character(),  
##   latDMS = col_character(),  
##   lonDMS = col_character(),  
##   textLatDD = col_character()
```

```
## )

rawData %>%
  select(catalogNumber, textLatDD)

## # A tibble: 10,767 x 2
##   catalogNumber textLatDD
##         <dbl> <chr>
## 1      219088 37.7609527778
## 2      233524 37.8999569
## 3      234346 37.8999569
## 4      233951 37.8999569
## 5      235290 37.8999569
## 6       85106 43.2751187
## 7       50048 34.53903
## 8      216309 36.97099
## 9      294933 36.584948
## 10     50255 44.2611111111
## # ... with 10,757 more rows

# test the creation of a numLatDD column as a numeric column and see what rows were converted to NA
rawData %>%
  mutate(numLatDD = as.numeric(rawData$textLatDD)) %>%
  filter(is.na(numLatDD)) %>%
  select(textLatDD, numLatDD) %>%
  print() %>%
  group_by(textLatDD) %>%
  summarize(count = n())

## Warning: NAs introduced by coercion

## # A tibble: 1,222 x 2
##   textLatDD numLatDD
##   <chr>      <dbl>
## 1 missing      NA
## 2 missing      NA
## 3 missing      NA
## 4 missing      NA
## 5 missing      NA
## 6 missing      NA
## 7 missing      NA
## 8 missing      NA
## 9 missing      NA
## 10 missing     NA
## # ... with 1,212 more rows

## # A tibble: 1 x 2
##   textLatDD count
##   <chr>      <int>
## 1 missing    1222

# create a numeric column based on the previously tested conversion of the textLatDD column
rawData$numLatDD <- as.numeric(rawData$textLatDD)

## Warning: NAs introduced by coercion

rawData
```

```
## # A tibble: 10,767 x 25
##       X1 uuid  institutionCode collectionCode catalogNumber recordedBy
##   <dbl> <chr> <chr>           <chr>           <dbl> <chr>
## 1     1 0603~ mvz             mammal specim~    219088 collector~
## 2     2 0fb1~ mvz             mammal specim~    233524 collector~
## 3     3 1a69~ mvz             mammal specim~    234346 collector~
## 4     4 1a99~ mvz             mammal specim~    233951 collector~
## 5     5 1f3b~ mvz             mammal specim~    235290 collector~
## 6     6 203f~ uam             mammal specim~     85106 collector~
## 7     7 21ce~ omnh            mammals           50048 caldwell,~
## 8     8 23d3~ mvz             mammal specim~    216309 collector~
## 9     9 244b~ msb             mammal specim~    294933 collector~
## 10    10 2682~ uam             mammal specim~     50255 collector~
## # ... with 10,757 more rows, and 19 more variables: countryCode <chr>,
## #   stateProvince <chr>, county <chr>, decimalLatitude <dbl>,
## #   decimalLongitude <dbl>, eventDate <dtm>, year <dbl>, month <dbl>,
## #   day <dbl>, genus <chr>, specificEpithet <chr>, scientificName <chr>,
## #   weight <dbl>, length <dbl>, sex <chr>, latDMS <chr>, lonDMS <chr>,
## #   textLatDD <chr>, numLatDD <dbl>
```

We can also accomplish a similar outcome by specifying the column type that should be created as part of the import process.

```
rawData2 <- read_csv("../data/learning.csv",
  col_types = cols(
    textLatDD = col_double()
  ),
  progress = FALSE)
```

```
## Warning: Missing column names filled in: 'X1' [1]
```

```
## Warning: 1222 parsing failures.
```

```
## row      col expected actual      file
## 11 textLatDD a double missing '../data/learning.csv'
## 21 textLatDD a double missing '../data/learning.csv'
## 25 textLatDD a double missing '../data/learning.csv'
## 32 textLatDD a double missing '../data/learning.csv'
## 39 textLatDD a double missing '../data/learning.csv'
## ... ..
## See problems(...) for more details.
```

```
# Display the column definitions for the imported dataset
spec(rawData2)
```

```
## cols(
##   X1 = col_double(),
##   uuid = col_character(),
##   institutionCode = col_character(),
##   collectionCode = col_character(),
##   catalogNumber = col_double(),
##   recordedBy = col_character(),
##   countryCode = col_character(),
##   stateProvince = col_character(),
##   county = col_character(),
##   decimalLatitude = col_double(),
##   decimalLongitude = col_double(),
```



```

##   eventDate = col_datetime(format = ""),
##   year = col_double(),
##   month = col_double(),
##   day = col_double(),
##   genus = col_character(),
##   specificEpithet = col_character(),
##   scientificName = col_character(),
##   weight = col_double(),
##   length = col_double(),
##   sex = col_character(),
##   latDMS = col_character(),
##   lonDMS = col_character(),
##   textLatDD = col_double()
## )

# Report the problems that were encountered when the data were imported.
problems(rawData2)

## # A tibble: 1,222 x 5
##   row col      expected actual  file
##   <int> <chr>      <chr>    <chr>  <chr>
## 1    11 textLatDD a double missing './data/learning.csv'
## 2    21 textLatDD a double missing './data/learning.csv'
## 3    25 textLatDD a double missing './data/learning.csv'
## 4    32 textLatDD a double missing './data/learning.csv'
## 5    39 textLatDD a double missing './data/learning.csv'
## 6    44 textLatDD a double missing './data/learning.csv'
## 7    54 textLatDD a double missing './data/learning.csv'
## 8    55 textLatDD a double missing './data/learning.csv'
## 9    59 textLatDD a double missing './data/learning.csv'
## 10   63 textLatDD a double missing './data/learning.csv'
## # ... with 1,212 more rows

# Display the imported table
rawData2

## # A tibble: 10,767 x 24
##   X1 uuid  institutionCode collectionCode catalogNumber recordedBy
##   <dbl> <chr> <chr>                <chr>          <dbl> <chr>
## 1     1 0603~ mvz                mammal specim~    219088 collector~
## 2     2 0fb1~ mvz                mammal specim~    233524 collector~
## 3     3 1a69~ mvz                mammal specim~    234346 collector~
## 4     4 1a99~ mvz                mammal specim~    233951 collector~
## 5     5 1f3b~ mvz                mammal specim~    235290 collector~
## 6     6 203f~ uam                mammal specim~     85106 collector~
## 7     7 21ce~ omnh                mammals                50048 caldwell,~
## 8     8 23d3~ mvz                mammal specim~    216309 collector~
## 9     9 244b~ msb                mammal specim~    294933 collector~
## 10   10 2682~ uam                mammal specim~     50255 collector~
## # ... with 10,757 more rows, and 18 more variables: countryCode <chr>,
## #   stateProvince <chr>, county <chr>, decimalLatitude <dbl>,
## #   decimalLongitude <dbl>, eventDate <dtm>, year <dbl>, month <dbl>,
## #   day <dbl>, genus <chr>, specificEpithet <chr>, scientificName <chr>,
## #   weight <dbl>, length <dbl>, sex <chr>, latDMS <chr>, lonDMS <chr>,
## #   textLatDD <dbl>

```

```
# convert the catalogNumberTxt column to a character column and see what the result is
rawData %>%
```

```
  mutate(catalogNumberTxt = as.character(catalogNumber)) %>%
  filter(is.na(catalogNumberTxt))
```

```
## # A tibble: 0 x 26
## #   ... with 26 variables: X1 <dbl>, uuid <chr>, institutionCode <chr>,
## #     collectionCode <chr>, catalogNumber <dbl>, recordedBy <chr>,
## #     countryCode <chr>, stateProvince <chr>, county <chr>,
## #     decimalLatitude <dbl>, decimalLongitude <dbl>, eventDate <dtm>,
## #     year <dbl>, month <dbl>, day <dbl>, genus <chr>, specificEpithet <chr>,
## #     scientificName <chr>, weight <dbl>, length <dbl>, sex <chr>, latDMS <chr>,
## #     lonDMS <chr>, textLatDD <chr>, numLatDD <dbl>, catalogNumberTxt <chr>
```

### 3. Check and handle missing values

It is important to understand the potential impact that missing data will have on your analysis. As we've already seen the import process may automatically produce missing data values in your analysis dataframe (or *tibble* in the context of tidyverse based processes). Some functions enable you to efficiently visualize the patterns of missing values in your dataset - allowing for the analysis of large datasets that otherwise would not be feasible to review manually.

```
paste("decimalLatitude: number of NA values", sum(is.na(rawData$decimalLatitude)), sep = " ")
```

```
## [1] "decimalLatitude: number of NA values 1222"
```

```
paste("decimalLongitude: number of NA values", sum(is.na(rawData$decimalLongitude)), sep = " ")
```

```
## [1] "decimalLongitude: number of NA values 1222"
```

```
paste("weight: number of NA values", sum(is.na(rawData$weight)), sep = " ")
```

```
## [1] "weight: number of NA values 0"
```

```
paste("length: number of NA values", sum(is.na(rawData$length)), sep = " ")
```

```
## [1] "length: number of NA values 0"
```

```
paste("sex: number of NA values", sum(is.na(rawData$sex)), sep = " ")
```

```
## [1] "sex: number of NA values 1446"
```

```
paste("latDMS: number of NA values", sum(is.na(rawData$latDMS)), sep = " ")
```

```
## [1] "latDMS: number of NA values 0"
```

```
paste("lonDMS: number of NA values", sum(is.na(rawData$lonDMS)), sep = " ")
```

```
## [1] "lonDMS: number of NA values 0"
```

```
paste("textLatDD: number of NA values", sum(is.na(rawData$textLatDD)), sep = " ")
```

```
## [1] "textLatDD: number of NA values 0"
```

```
paste("numLatDD: number of NA values", sum(is.na(rawData$numLatDD)), sep = " ")
```

```
## [1] "numLatDD: number of NA values 1222"
```

In this analysis we will be using the `md.pattern` function that is part of the `mice` and `VIM` packages. If you haven't already installed the `mice` package you can do so by executing the `install.packages("mice")` command.

```
##
## Attaching package: 'mice'

## The following objects are masked from 'package:base':
##
##      cbind, rbind

## Loading required package: colorspace
## Loading required package: grid
## Loading required package: data.table

##
## Attaching package: 'data.table'

## The following objects are masked from 'package:dplyr':
##
##      between, first, last

## The following object is masked from 'package:purrr':
##
##      transpose

## VIM is ready to use.
## Since version 4.0.0 the GUI is in its own package VIMGUI.
##
##      Please use the package to use the new (and old) GUI.

## Suggestions and bug-reports can be submitted at: https://github.com/alexkowa/VIM/issues

##
## Attaching package: 'VIM'

## The following object is masked from 'package:datasets':
##
##      sleep
```

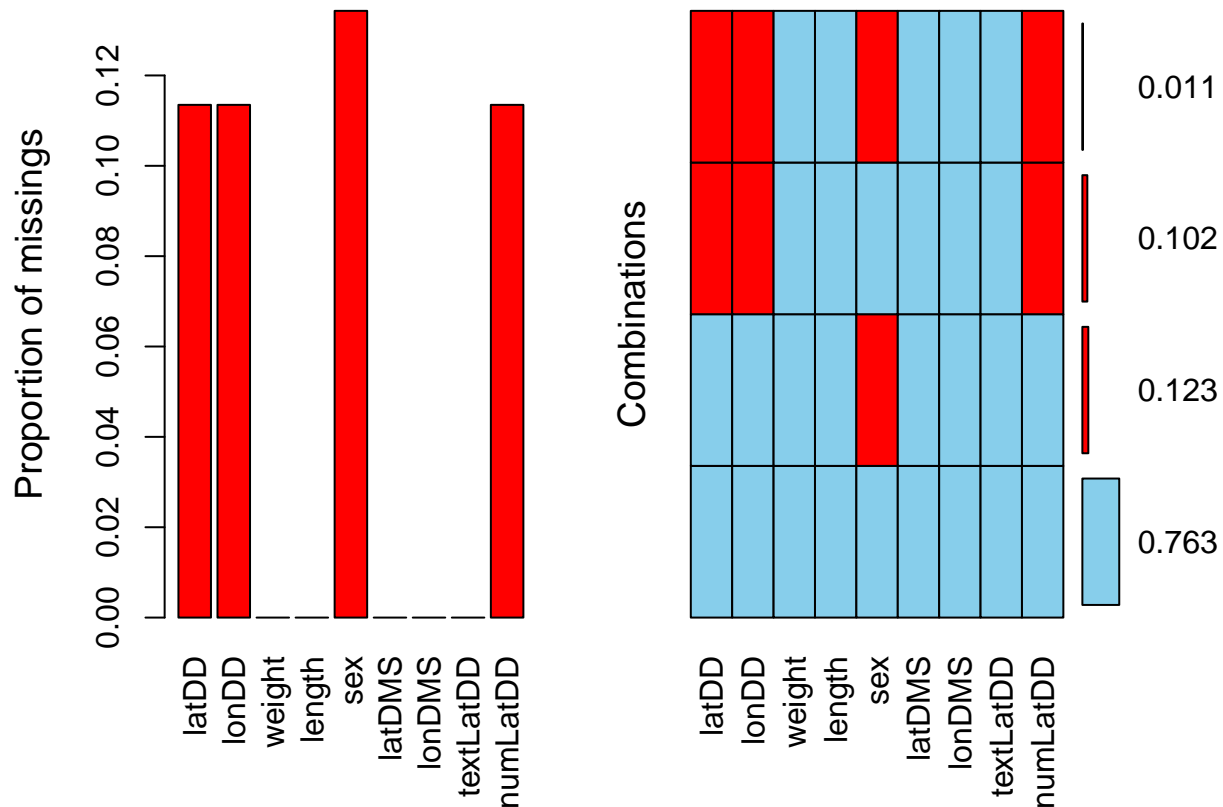
	weight	length	latDMS	lonDMS	textLatDD	decimalLatitude	decimalLongitude	numLatDD	sex	
8218										0
1327										1
1103										3
119										4
	0	0	0	0	0	12221	2221	2221	1445	112

```
##      weight length latDMS lonDMS textLatDD decimalLatitude decimalLongitude numLatDD sex
## 8218      1      1      1      1      1      1      1      1      1      1      0
```

```
## 1327      1      1      1      1      1      1      1      1      0      1
## 1103      1      1      1      1      1      0      0      0      1      3
## 119       1      1      1      1      1      0      0      0      0      4
##          0      0      0      0      0      1222    1222    1222 1446 5112
```

Another method of viewing similar information

```
rawData %>%
  select(decimalLatitude, decimalLongitude, weight, length, sex, latDMS, lonDMS, textLatDD, numLatDD) %>%
  rename(latDD = decimalLatitude, lonDD = decimalLongitude) %>%
  aggr(numbers=TRUE)
```



#### 4. Multi-value columns

A core principle of having well structured data that is read for analysis is that:

In the context of “Tidy” data that underlie the tools developed as part of the tidyverse package<sup>2</sup>[Hadley Wickham & Garrett Golemund. 2017. *R for Data Science*. O’Reilly. - section on Tidy Data <https://r4ds.had.co.nz/tidy-data.html>

There are three interrelated rules which make a dataset tidy:

1. Each variable must have its own column.
2. Each observation must have its own row.
3. **Each value must have its own cell.**

This issue also relates to the idea of *atomicity* in Codd’s definition of *First Normal Form* when *normalizing* a relational database<sup>2</sup>. While we’re not going to get into relational data modeling in R here, well structured data allow for the use of relational data models in your analysis independent of a separate database server.

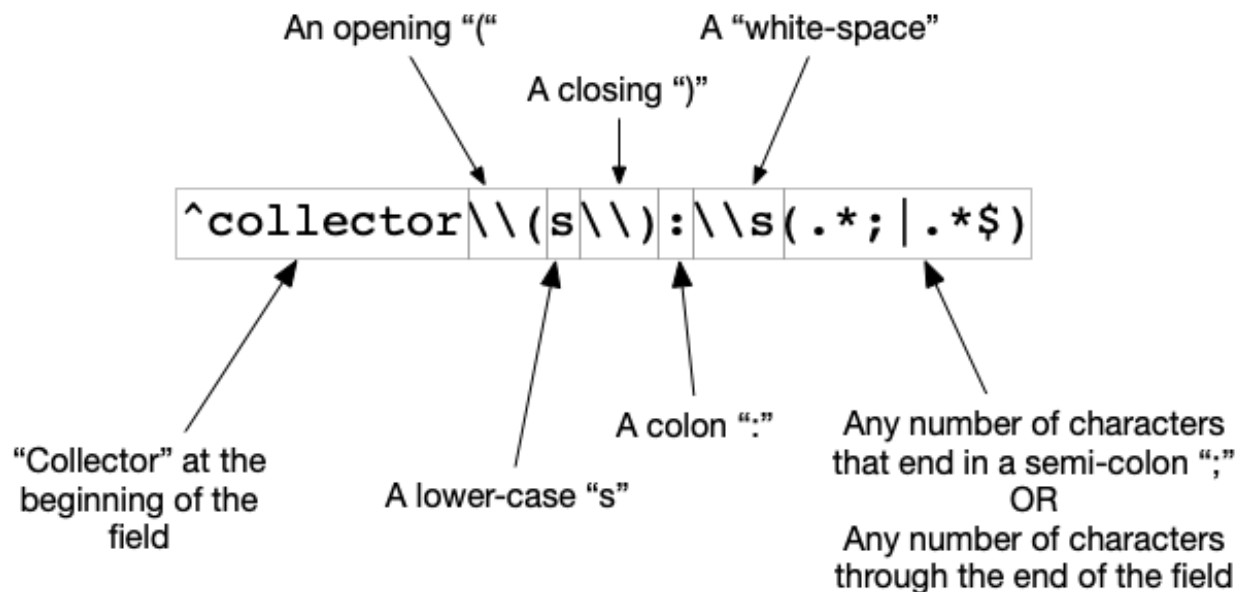
<sup>2</sup>[https://en.wikipedia.org/wiki/First\\_normal\\_form](https://en.wikipedia.org/wiki/First_normal_form)

In this example we are going to focus on three columns: `recordedBy`, `latDMS`, and `lonDMS`.

```
rawData %>%
  select(recordedBy, latDMS, lonDMS)
```

```
## # A tibble: 10,767 x 3
##   recordedBy                                latDMS                lonDMS
##   <chr>                                     <chr>                 <chr>
## 1 collector(s): ana lilia trujano álvarez, eric ghilarducci "37°45'39.430\"N" "-122° 6'48.370\"E"
## 2 collector(s): william z. lidicker jr.                    "37°53'59.845\"N" "-123°38'18.039\"E"
## 3 collector(s): william z. lidicker jr.                    "37°53'59.845\"N" "-123°38'18.039\"E"
## 4 collector(s): william z. lidicker jr.                    "37°53'59.845\"N" "-123°38'18.039\"E"
## 5 collector(s): william z. lidicker jr.                    "37°53'59.845\"N" "-123°38'18.039\"E"
## 6 collector(s): tom manning; preparator(s): amber baxter   "43°16'30.427\"N" "-123°12'32.023\"E"
## 7 caldwell, j. p. and vitt, l. j.                          "34°32'20.508\"N" "-95° 6'42.444\"E"
## 8 collector(s): james l. patton                             "36°58'15.564\"N" "-119°54'16.740\"E"
## 9 collector(s): troy l. best; preparator(s): troy l. best  "36°35' 5.813\"N" "-108° 0'55.757\"E"
## 10 collector(s): karl j. martin; preparator(s): paul ollig "44°15'40.000\"N" "-124°25' 1.000\"E"
## # ... with 10,757 more rows
```

Breaking apart the `recordedBy` column using the `str_match` function from the `stringr` package. This uses *regular expressions* for defining the text patterns that should be found and processed in the process of breaking the column apart into new columns. Regular expressions are an art to themselves and there are many resources for learning their effective use - ranging from one-page cheat-sheets to full length books. The following figure describes the structure of the regular expressions used in the sample code:



**Note:** In R standard RegEx escaped characters (like `\"` for `(`, `\"` for `)`, and `\\s` for white-space) need to have an additional backslash (`\"`) added before them for R to properly translate the RegEx syntax

Figure 2: Description of the regular expression used to extract targeted substrings from the combined `recordedBy` field

```

collectorExtract <- "^collector\\(s\\):\\s(.*;|.*)$"
preparatorExtract <- "preparator\\(s\\):\\s(.*;|.*)$"
#str_match(rawData$recordedBy, collectorExtract)[,2]
#str_match(rawData$recordedBy, preparatorExtract)[,2]
rawData$collectors <- str_match(rawData$recordedBy, collectorExtract)[,2]
rawData$preparators <- str_match(rawData$recordedBy, preparatorExtract)[,2]

```

What would the next logical step in the process be for cleaning up the *recordedBy* column?

Breaking apart the latDMS and lonDMS columns into their constituent parts

```

dmsExtract <- "\\s*(-*[:digit:]+)°\\s*([[:digit:]+)\\'\\s*([[:digit:]+)"

latSubstrings <- str_match(rawData$latDMS, dmsExtract)
rawData$latD <- as.numeric(latSubstrings[,2])
rawData$latM <- as.numeric(latSubstrings[,3])
rawData$latS <- as.numeric(latSubstrings[,4])

glimpse(latSubstrings)

```

```
## chr [1:10767, 1:4] "37°45'39" "37°53'59" "37°53'59" "37°53'59" "37°53'59" "43°16'30" "34°32'20" "36°
```

```

lonSubstrings <- str_match(rawData$lonDMS, dmsExtract)
rawData$lonD <- as.numeric(lonSubstrings[,2])
rawData$lonM <- as.numeric(lonSubstrings[,3])
rawData$lonS <- as.numeric(lonSubstrings[,4])

glimpse(lonSubstrings)

```

```
## chr [1:10767, 1:4] "-122° 6'48" "-123°38'18" "-123°38'18" "-123°38'18" "-123°38'18" "-123°12'32" "-123°12'32"
```

```
glimpse(rawData)
```

```

## Observations: 10,767
## Variables: 33
## $ X1 <dbl> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000, 1001, 1002, 1003, 1004, 1005, 1006, 1007, 1008, 1009, 1010, 1011, 1012, 1013, 1014, 1015, 1016, 1017, 1018, 1019, 1020, 1021, 1022, 1023, 1024, 1025, 1026, 1027, 1028, 1029, 1030, 1031, 1032, 1033, 1034, 1035, 1036, 1037, 1038, 1039, 1040, 1041, 1042, 1043, 1044, 1045, 1046, 1047, 1048, 1049, 1050, 1051, 1052, 1053, 1054, 1055, 1056, 1057, 1058, 1059, 1060, 1061, 1062, 1063, 1064, 1065, 1066, 1067, 1068, 1069, 1070, 1071, 1072, 1073, 1074, 1075, 1076, 1077, 1078, 1079, 1080, 1081, 1082, 1083, 1084, 1085, 1086, 1087, 1088, 1089, 1090, 1091, 1092, 1093, 1094, 1095, 1096, 1097, 1098, 1099, 1100, 1101, 1102, 1103, 1104, 1105, 1106, 1107, 1108, 1109, 1110, 1111, 1112, 1113, 1114, 1115, 1116, 1117, 1118, 1119, 1120, 1121, 1122, 1123, 1124, 1125, 1126, 1127, 1128, 1129, 1130, 1131, 1132, 1133, 1134, 1135, 1136, 1137, 1138, 1139, 1140, 1141, 1142, 1143, 1144, 1145, 1146, 1147, 1148, 1149, 1150, 1151, 1152, 1153, 1154, 1155, 1156, 1157, 1158, 1159, 1160, 1161, 1162, 1163, 1164, 1165, 1166, 1167, 1168, 1169, 1170, 1171, 1172, 1173, 1174, 1175, 1176, 1177, 1178, 1179, 1180, 1181, 1182, 1183, 1184, 1185, 1186, 1187, 1188, 1189, 1190, 1191, 1192, 1193, 1194, 1195, 1196, 1197, 1198, 1199, 1200, 1201, 1202, 1203, 1204, 1205, 1206, 1207, 1208, 1209, 1210, 1211, 1212, 1213, 1214, 1215, 1216, 1217, 1218, 1219, 1220, 1221, 1222, 1223, 1224, 1225, 1226, 1227, 1228, 1229, 1230, 1231, 1232, 1233, 1234, 1235, 1236, 1237, 1238, 1239, 1240, 1241, 1242, 1243, 1244, 1245, 1246, 1247, 1248, 1249, 1250, 1251, 1252, 1253, 1254, 1255, 1256, 1257, 1258, 1259, 1260, 1261, 1262, 1263, 1264, 1265, 1266, 1267, 1268, 1269, 1270, 1271, 1272, 1273, 1274, 1275, 1276, 1277, 1278, 1279, 1280, 1281, 1282, 1283, 1284, 1285, 1286, 1287, 1288, 1289, 1290, 1291, 1292, 1293, 1294, 1295, 1296, 1297, 1298, 1299, 1300, 1301, 1302, 1303, 1304, 1305, 1306, 1307, 1308, 1309, 1310, 1311, 1312, 1313, 1314, 1315, 1316, 1317, 1318, 1319, 1320, 1321, 1322, 1323, 1324, 1325, 1326, 1327, 1328, 1329, 1330, 1331, 1332, 1333, 1334, 1335, 1336, 1337, 1338, 1339, 1340, 1341, 1342, 1343, 1344, 1345, 1346, 1347, 1348, 1349, 1350, 1351, 1352, 1353, 1354, 1355, 1356, 1357, 1358, 1359, 1360, 1361, 1362, 1363, 1364, 1365, 1366, 1367, 1368, 1369, 1370, 1371, 1372, 1373, 1374, 1375, 1376, 1377, 1378, 1379, 1380, 1381, 1382, 1383, 1384, 1385, 1386, 1387, 1388, 1389, 1390, 1391, 1392, 1393, 1394, 1395, 1396, 1397, 1398, 1399, 1400, 1401, 1402, 1403, 1404, 1405, 1406, 1407, 1408, 1409, 1410, 1411, 1412, 1413, 1414, 1415, 1416, 1417, 1418, 1419, 1420, 1421, 1422, 1423, 1424, 1425, 1426, 1427, 1428, 1429, 1430, 1431, 1432, 1433, 1434, 1435, 1436, 1437, 1438, 1439, 1440, 1441, 1442, 1443, 1444, 1445, 1446, 1447, 1448, 1449, 1450, 1451, 1452, 1453, 1454, 1455, 1456, 1457, 1458, 1459, 1460, 1461, 1462, 1463, 1464, 1465, 1466, 1467, 1468, 1469, 1470, 1471, 1472, 1473, 1474, 1475, 1476, 1477, 1478, 1479, 1480, 1481, 1482, 1483, 1484, 1485, 1486, 1487, 1488, 1489, 1490, 1491, 1492, 1493, 1494, 1495, 1496, 1497, 1498, 1499, 1500, 1501, 1502, 1503, 1504, 1505, 1506, 1507, 1508, 1509, 1510, 1511, 1512, 1513, 1514, 1515, 1516, 1517, 1518, 1519, 1520, 1521, 1522, 1523, 1524, 1525, 1526, 1527, 1528, 1529, 1530, 1531, 1532, 1533, 1534, 1535, 1536, 1537, 1538, 1539, 1540, 1541, 1542, 1543, 1544, 1545, 1546, 1547, 1548, 1549, 1550, 1551, 1552, 1553, 1554, 1555, 1556, 1557, 1558, 1559, 1560, 1561, 1562, 1563, 1564, 1565, 1566, 1567, 1568, 1569, 1570, 1571, 1572, 1573, 1574, 1575, 1576, 1577, 1578, 1579, 1580, 1581, 1582, 1583, 1584, 1585, 1586, 1587, 1588, 1589, 1590, 1591, 1592, 1593, 1594, 1595, 1596, 1597, 1598, 1599, 1600, 1601, 1602, 1603, 1604, 1605, 1606, 1607, 1608, 1609, 1610, 1611, 1612, 1613, 1614, 1615, 1616, 1617, 1618, 1619, 1620, 1621, 1622, 1623, 1624, 1625, 1626, 1627, 1628, 1629, 1630, 1631, 1632, 1633, 1634, 1635, 1636, 1637, 1638, 1639, 1640, 1641, 1642, 1643, 1644, 1645, 1646, 1647, 1648, 1649, 1650, 1651, 1652, 1653, 1654, 1655, 1656, 1657, 1658, 1659, 1660, 1661, 1662, 1663, 1664, 1665, 1666, 1667, 1668, 1669, 1670, 1671, 1672, 1673, 1674, 1675, 1676, 1677, 1678, 1679, 1680, 1681, 1682, 1683, 1684, 1685, 1686, 1687, 1688, 1689, 1690, 1691, 1692, 1693, 1694, 1695, 1696, 1697, 1698, 1699, 1700, 1701, 1702, 1703, 1704, 1705, 1706, 1707, 1708, 1709, 1710, 1711, 1712, 1713, 1714, 1715, 1716, 1717, 1718, 1719, 1720, 1721, 1722, 1723, 1724, 1725, 1726, 1727, 1728, 1729, 1730, 1731, 1732, 1733, 1734, 1735, 1736, 1737, 1738, 1739, 1740, 1741, 1742, 1743, 1744, 1745, 1746, 1747, 1748, 1749, 1750, 1751, 1752, 1753, 1754, 1755, 1756, 1757, 1758, 1759, 1760, 1761, 1762, 1763, 1764, 1765, 1766, 1767, 1768, 1769, 1770, 1771, 1772, 1773, 1774, 1775, 1776, 1777, 1778, 1779, 1780, 1781, 1782, 1783, 1784, 1785, 1786, 1787, 1788, 1789, 1790, 1791, 1792, 1793, 1794, 1795, 1796, 1797, 1798, 1799, 1800, 1801, 1802, 1803, 1804, 1805, 1806, 1807, 1808, 1809, 1810, 1811, 1812, 1813, 1814, 1815, 1816, 1817, 1818, 1819, 1820, 1821, 1822, 1823, 1824, 1825, 1826, 1827, 1828, 1829, 1830, 1831, 1832, 1833, 1834, 1835, 1836, 1837, 1838, 1839, 1840, 1841, 1842, 1843, 1844, 1845, 1846, 1847, 1848, 1849, 1850, 1851, 1852, 1853, 1854, 1855, 1856, 1857, 1858, 1859, 1860, 1861, 1862, 1863, 1864, 1865, 1866, 1867, 1868, 1869, 1870, 1871, 1872, 1873, 1874, 1875, 1876, 1877, 1878, 1879, 1880, 1881, 1882, 1883, 1884, 1885, 1886, 1887, 1888, 1889, 1890, 1891, 1892, 1893, 1894, 1895, 1896, 1897, 1898, 1899, 1900, 1901, 1902, 1903, 1904, 1905, 1906, 1907, 1908, 1909, 1910, 1911, 1912, 1913, 1914, 1915, 1916, 1917, 1918, 1919, 1920, 1921, 1922, 1923, 1924, 1925, 1926, 1927, 1928, 1929, 1930, 1931, 1932, 1933, 1934, 1935, 1936, 1937, 1938, 1939, 1940, 1941, 1942, 1943, 1944, 1945, 1946, 1947, 1948, 1949, 1950, 1951, 1952, 1953, 1954, 1955, 1956, 1957, 1958, 1959, 1960, 1961, 1962, 1963, 1964, 1965, 1966, 1967, 1968, 1969, 1970, 1971, 1972, 1973, 1974, 1975, 1976, 1977, 1978, 1979, 1980, 1981, 1982, 1983, 1984, 1985, 1986, 1987, 1988, 1989, 1990, 1991, 1992, 1993, 1994, 1995, 1996, 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030, 2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043, 2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056, 2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069, 2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082, 2083, 2084, 2085, 2086, 2087, 2088, 2089, 2090, 2091, 2092, 2093, 2094, 2095, 2096, 20
```

```
## $ latDMS      <chr> "37°45'39.430\"N", "37°53'59.845\"N", "37°53'59.845\"N", "37°53'59.845\"N",
## $ lonDMS      <chr> "-122° 6'48.370\"E", "-123°38'18.039\"E", "-123°38'18.039\"E", "-123°38'18.039\"E",
## $ textLatDD    <chr> "37.7609527778", "37.8999569", "37.8999569", "37.8999569", "37.8999569", "43.27512",
## $ numLatDD     <dbl> 37.76095, 37.89996, 37.89996, 37.89996, 37.89996, 43.27512, 34.53903, 36.97000,
## $ collectors   <chr> "ana lilia trujano álvarez, eric ghilarducci", "william z. lidicker jr.", "troy l. best",
## $ preparators  <chr> NA, NA, NA, NA, NA, "amber baxter", NA, NA, "troy l. best", "paul ollig", NA, NA, NA, NA, NA, NA,
## $ latD         <dbl> 37, 37, 37, 37, 37, 43, 34, 36, 36, 44, NA, 39, 45, 37, 37, 32, 37, 36, 36,
## $ latM         <dbl> 45, 53, 53, 53, 53, 16, 32, 58, 35, 15, NA, 59, 32, 53, 30, 30, 53, 28, 58,
## $ latS         <dbl> 39, 59, 59, 59, 59, 30, 20, 15, 5, 40, NA, 36, 24, 59, 12, 27, 59, 39, 1, 1,
## $ lonD         <dbl> -122, -123, -123, -123, -123, -123, -123, -95, -119, -108, -124, NA, -109, -109, -109, -109, -109, -109,
## $ lonM         <dbl> 6, 38, 38, 38, 38, 12, 6, 54, 0, 25, NA, 46, 10, 38, 51, 17, 38, 46, 41, 32,
## $ lonS         <dbl> 48, 18, 18, 18, 18, 32, 42, 16, 55, 1, NA, 28, 48, 18, 8, 58, 18, 52, 6, 6,
```

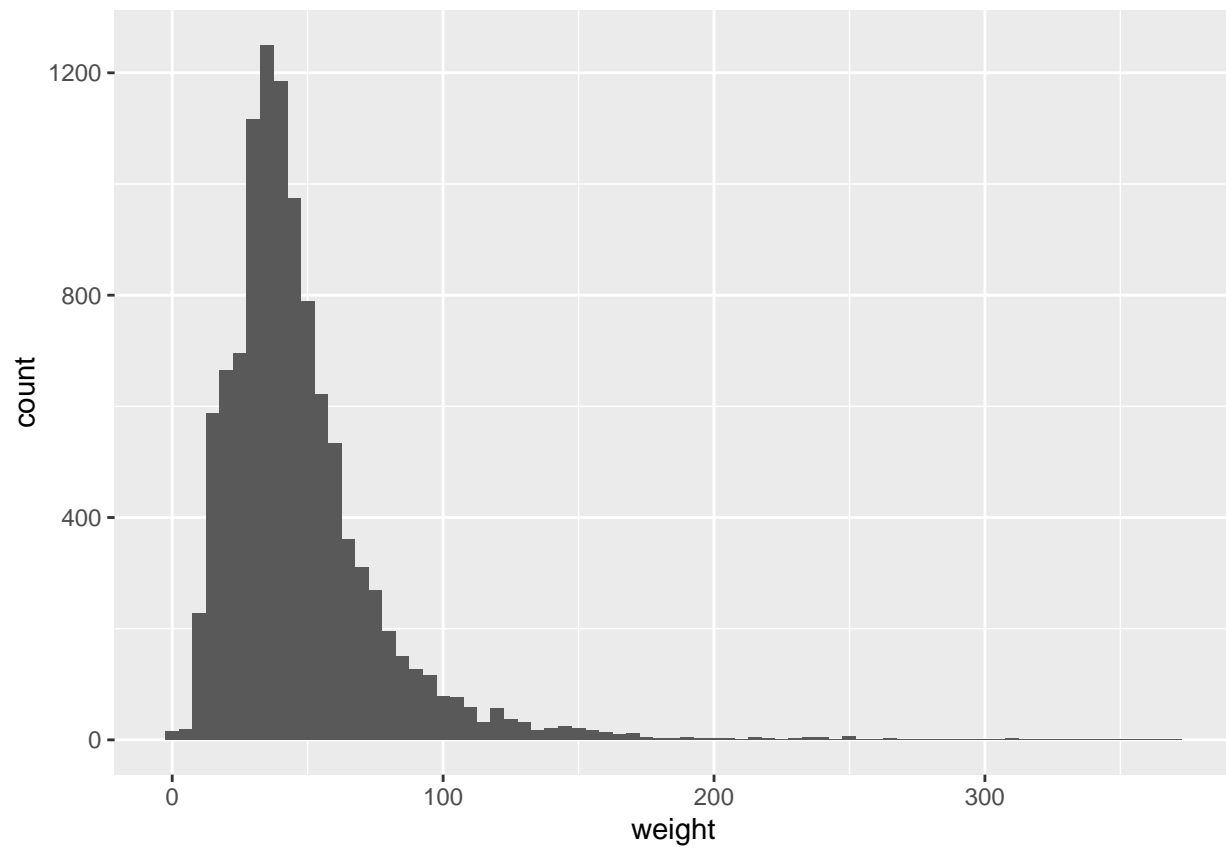
## 5. Check value ranges and explore data

As part of the examination of these columns we will use the `assertr`. If you need to install the `assertr` package in your environment you can execute the `install.packages("assertr")` command.

Checking the `weight` and `length` columns.

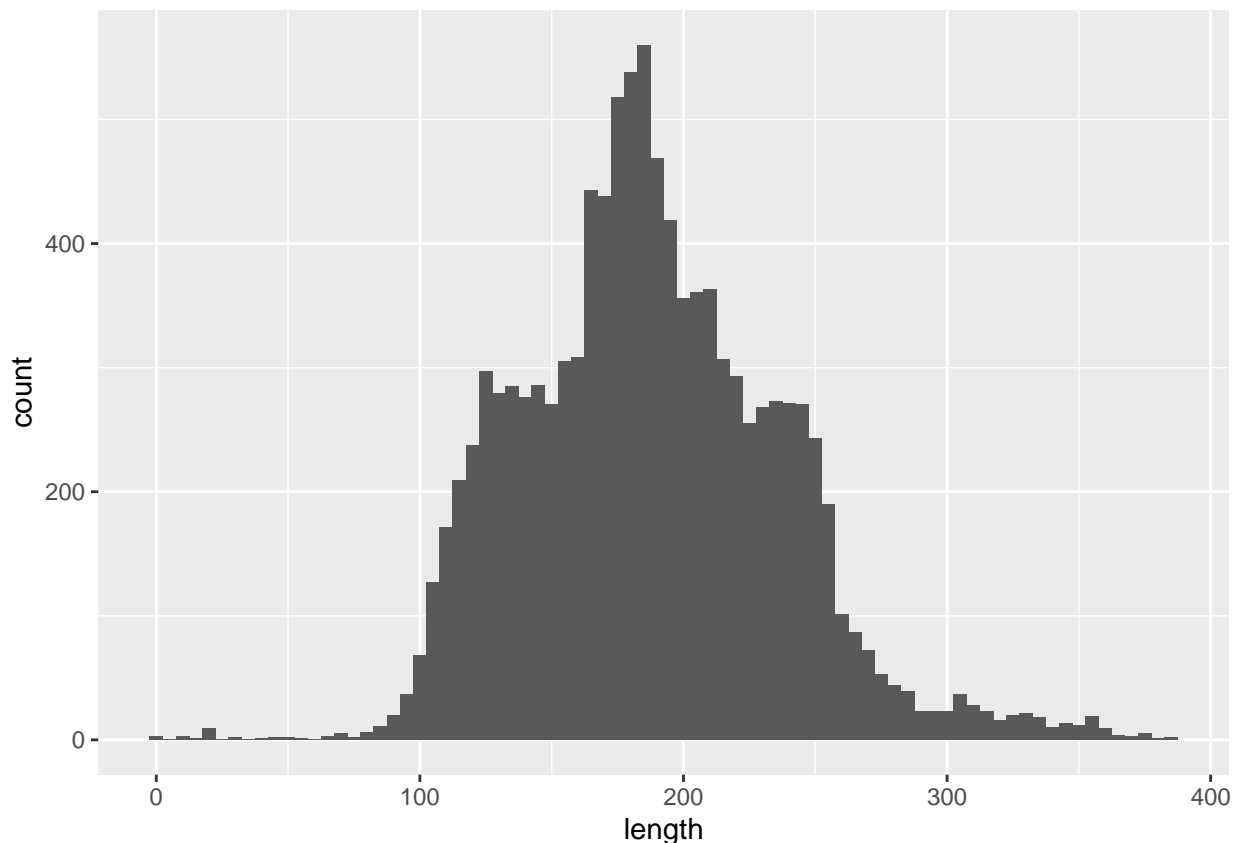
```
#library(assertr)
#
#rawData %>%
#  chain_start %>%
#  assert(within_bounds(1,Inf), weight) %>%
#  assert(within_bounds(1,Inf), length) %>%
#  insist(within_n_sds(3), weight) %>%
#  insist(within_n_sds(3), length) %>%
#  chain_end
```

```
library(ggplot2)
ggplot(rawData, aes(x=weight)) +
  geom_histogram(binwidth = 5)
```



```
ggplot(rawData, aes(x=length)) +  
  geom_histogram(binwidth = 5)
```





## 6. Bringing it all together to produce our analysis dataset

```
# import the data with explicit column definitions for the textLatDD and catalogNumber columns
analysisData <- read_csv("../data/learning.csv",
  col_types = cols(
    textLatDD = col_double(),
    catalogNumber = col_character()
  ),
  progress = FALSE)
```

```
## Warning: Missing column names filled in: 'X1' [1]
```

```
## Warning: 1222 parsing failures.
```

```
## row      col expected actual      file
## 11 textLatDD a double missing '../data/learning.csv'
## 21 textLatDD a double missing '../data/learning.csv'
## 25 textLatDD a double missing '../data/learning.csv'
## 32 textLatDD a double missing '../data/learning.csv'
## 39 textLatDD a double missing '../data/learning.csv'
## ... .....
## See problems(...) for more details.
```

```
# split up the recordedBy column
```

```
collectorExtract <- "^collector\\(s\\):\\s(.+;|.*)$"
```

```
preparatorExtract <- "preparator\\(s\\):\\s(.+;|.*)$"
```

```
analysisData$collectors <- str_match(analysisData$recordedBy, collectorExtract)[,2]
```

```
analysisData$preparators <- str_match(analysisData$recordedBy, preparatorExtract)[,2]
```

```

# split up the latDMS and lonDMS columns
dmsExtract <- "\\s*(-*[:digit:]+)°\\s*([:digit:]+)\\'\\s*([:digit:]+)"

latSubstrings <- str_match(analysisData$latDMS, dmsExtract)
analysisData$latD <- as.numeric(latSubstrings[,2])
analysisData$latM <- as.numeric(latSubstrings[,3])
analysisData$latS <- as.numeric(latSubstrings[,4])

lonSubstrings <- str_match(analysisData$lonDMS, dmsExtract)
analysisData$lonD <- as.numeric(lonSubstrings[,2])
analysisData$lonM <- as.numeric(lonSubstrings[,3])
analysisData$lonS <- as.numeric(lonSubstrings[,4])

glimpse(analysisData)

```

```

## Observations: 10,767
## Variables: 32
## $ X1 <dbl> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, ...
## $ uuid <chr> "060380ea-7b06-474e-8d2e-b6e4a8c21e1a", "0fb17a79-a8ce-45b6-b57a-2f640e8cccl
## $ institutionCode <chr> "mvz", "mvz", "mvz", "mvz", "mvz", "uam", "omnh", "mvz", "msb", "uam", "msb
## $ collectionCode <chr> "mammal specimens", "mammal specimens", "mammal specimens", "mammal specime
## $ catalogNumber <chr> "219088", "233524", "234346", "233951", "235290", "85106", "50048", "216309
## $ recordedBy <chr> "collector(s): ana lilia trujano álvarez, eric ghilarducci", "collector(s):
## $ countryCode <chr> "usa", "usa", "usa", "usa", "usa", "usa", "usa", "usa", "usa", "usa", "usa"
## $ stateProvince <chr> "california", "california", "california", "california", "california", "oreg
## $ county <chr> "contra costa county", "contra costa county", "contra costa county", "contr
## $ decimalLatitude <dbl> 37.76095, 37.89996, 37.89996, 37.89996, 37.89996, 43.27512, 34.53903, 36.97
## $ decimalLongitude <dbl> -121.88656, -122.36166, -122.36166, -122.36166, -122.36166, -122.79110, -94
## $ eventDate <dtm> 2005-11-23, 1959-06-21, 1962-11-22, 1960-07-31, 1964-07-04, 1996-10-23, 20
## $ year <dbl> 2005, 1959, 1962, 1960, 1964, 1996, 2011, 2005, 1989, 1996, 2013, 1977, 199
## $ month <dbl> 11, 6, 11, 7, 7, 10, 1, 8, 6, 5, 8, 8, 11, 6, 7, 6, 4, 8, 12, 10, 1, 9, 7, 8
## $ day <dbl> 22, 20, 21, 30, 3, 22, 17, 6, 4, 19, 10, 12, 30, 17, 17, 2, 21, 6, 26, 26, 1
## $ genus <chr> "microtus", "microtus", "microtus", "microtus", "microtus", "myodes", "micro
## $ specificEpithet <chr> "californicus", "californicus", "californicus", "californicus", "californi
## $ scientificName <chr> "microtus californicus californicus", "microtus californicus californicus",
## $ weight <dbl> 30.5, 22.0, 49.0, 33.0, 29.0, 23.5, 24.0, 27.0, 125.0, 12.0, 56.0, 33.5, 33
## $ length <dbl> 165, 143, 187, 169, 159, 141, 121, 176, 294, 110, 210, 148, 149, 170, 182, 1
## $ sex <chr> "male", "female", "male", "female", "female", "female", NA, "female", "fema
## $ latDMS <chr> "37°45'39.430\"N", "37°53'59.845\"N", "37°53'59.845\"N", "37°53'59.845\"N",
## $ lonDMS <chr> "-122° 6'48.370\"E", "-123°38'18.039\"E", "-123°38'18.039\"E", "-123°38'18.
## $ textLatDD <dbl> 37.76095, 37.89996, 37.89996, 37.89996, 37.89996, 43.27512, 34.53903, 36.97
## $ collectors <chr> "ana lilia trujano álvarez, eric ghilarducci", "william z. lidicker jr.", "
## $ preparators <chr> NA, NA, NA, NA, NA, "amber baxter", NA, NA, "troy l. best", "paul ollig", NA
## $ latD <dbl> 37, 37, 37, 37, 37, 43, 34, 36, 36, 44, NA, 39, 45, 37, 37, 32, 37, 36, 36,
## $ latM <dbl> 45, 53, 53, 53, 53, 16, 32, 58, 35, 15, NA, 59, 32, 53, 30, 30, 53, 28, 58,
## $ latS <dbl> 39, 59, 59, 59, 59, 30, 20, 15, 5, 40, NA, 36, 24, 59, 12, 27, 59, 39, 1, 1
## $ lonD <dbl> -122, -123, -123, -123, -123, -123, -95, -119, -108, -124, NA, -109, -109,
## $ lonM <dbl> 6, 38, 38, 38, 38, 12, 6, 54, 0, 25, NA, 46, 10, 38, 51, 17, 38, 46, 41, 32
## $ lonS <dbl> 48, 18, 18, 18, 18, 32, 42, 16, 55, 1, NA, 28, 48, 18, 8, 58, 18, 52, 6, 6,

```

Export the code from the workshop both as a documented and undocumented R script

```

#library(knitr)
#purl("cleaning_data.Rmd", "cleaning_data_nodocs.R", documentation = 0)

```