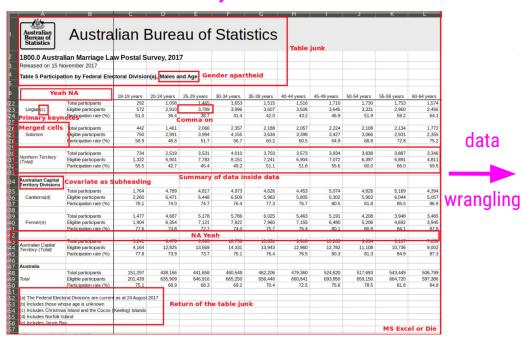
Tidying Data

Data Tidying

untidy data



tidy data

1	area	gender	age	State	Area (sq km)	Eligible participants	Participation rate (%)	Total participants	Total Paticipants
2	Adelaide	Female	18-19 years	SA	76	1341	83.5	1120	1120
3	Adelaide	Female	20-24 years	SA	76	4620	81.2	3750	3750
4	Adelaide	Female	25-29 years	SA	76	4897	81.8	4004	4004
5	Adelaide	Female	30-34 years	SA	76	4784	79.8	3820	3820
6	Adelaide	Female	35-39 years	SA	76	4319	79	3411	3411
7	Adelaide	Female	40-44 years	SA	76	4310	80.6	3472	3472
8	Adelaide	Female	45-49 years	SA	76	4579	81.4	3728	3728
9	Adelaide	Female	50-54 years	SA	76	4475	84.7	3791	3791
10	Adelaide	Female	55-59 years	SA	76	4622	87.3	4033	4033
11	Adelaide	Female	60-64 years	SA	76	4342	89.3	3879	3879
12	Adelaide	Female	65-69 years	SA	76	3970	90.7	3602	3602
13	Adelaide	Female	70-74 years	SA	76	3009	90.3	2716	2716
14	Adelaide	Female	75-79 years	SA	76	2156	88.5	1908	1908
15	Adelaide	Female	80-84 years	SA	76	1673	85.1	1423	1423

data

dplyr functions

```
# to install and load dplyr
install.packages("dplyr")
library(dplyr)
```

- %>% pipe operator for chaining a sequence of operations
- glimpse() get an overview of what's included in dataset
- filter() filter rows
- select() select, rename, and re-order columns
- rename() rename columns
- arrange() reorder rows
- mutate() create a new column
- group_by() group variables
- summarize() summarize information within a dataset
- left_join() combining data across data frame

tidyr functions

```
# to install and load tidyr
install.packages("tidyr")
library(tidyr)
```

- unite() combine contents of two or more columns into a single column
- separate() separate contents of a column into two or more columns

janitor functions

```
# to install and load janitor
install.packages("janitor")
library(janitor)
```

- clean_names() clean names of a data frame
- tabyl() get a helpful summary of a variable

skimr functions

```
# to install and load skimr
install.packages("skimr")
library(skimr)
```

• skim() - summarize a data frame

The pipe operator: %>%

```
If you want to: A \longrightarrow B
                                                                  Data frame A
                                                                  Function B()
    In R:
                                                                  Function C()
         Without the pipe operator B(A)
                                                                  Function D()
         With the pipe operator A %>% B
If you want to: A \longrightarrow B \longrightarrow C \longrightarrow D
    In R:
         Without the pipe operator D(C(B(A)))
         With the pipe operator A %>% B %>% C %>% D
```

https://dplyr.tidyverse.org/

Filtering Data —o—

```
> glimpse(msleep)
Observations: 83
```

\$ bodywt

```
Variables: 11
$ name
               <chr> "Cheetah", "Owl monkey", "Mountain beaver", "Greater short-tailed shrew", "Cow", "Three...
$ genus
               <chr> "Acinonyx", "Aotus", "Aplodontia", "Blarina", "Bos", "Bradypus", "Callorhinus", "Calomy...
               <chr> "carni", "omni", "herbi", "omni", "herbi", "herbi", "carni", NA, "carni", "herbi", "her...
$ vore
               <chr> "Carnivora", "Primates", "Rodentia", "Soricomorpha", "Artiodactyla", "Pilosa", "Carnivo...
$ order
$ conservation <chr> "lc", NA, "nt", "lc", "domesticated", NA, "vu", NA, "domesticated", "lc", "lc", "domest...
$ sleep_total
               <dbl> 12.1, 17.0, 14.4, 14.9, 4.0, 14.4, 8.7, 7.0, 10.1, 3.0, 5.3, 9.4, 10.0, 12.5, 10.3, 8.3...
$ sleep_rem
               <dbl> NA, 1.8, 2.4, 2.3, 0.7, 2.2, 1.4, NA, 2.9, NA, 0.6, 0.8, 0.7, 1.5, 2.2, 2.0, 1.4, 3.1, ...
$ sleep_cycle
               <dbl> NA, NA, NA, 0.133, 0.667, 0.767, 0.383, NA, 0.333, NA, NA, 0.217, NA, 0.117, NA, NA, 0....
$ awake
               <dbl> 11.9, 7.0, 9.6, 9.1, 20.0, 9.6, 15.3, 17.0, 13.9, 21.0, 18.7, 14.6, 14.0, 11.5, 13.7, 1...
$ brainwt
               <dbl> NA, 0.01550, NA, 0.00029, 0.42300, NA, NA, NA, 0.07000, 0.09820, 0.11500, 0.00550, NA, ...
```

<dbl> 50.000, 0.480, 1.350, 0.019, 600.000, 3.850, 20.490, 0.045, 14.000, 14.800, 33.500, 0.7...

```
> alimpse(msleep
Observations: 83
                     11 columns
Variables: 11
$ name
               <chr> "Cheetah", "Owl monkey", "Mountain beaver", "Greater short-tailed shrew", "Cow", "Three...
               <chr> "Acinonyx", "Aotus", "Aplodontia", "Blarina", "Bos", "Bradypus", "Callorhinus", "Calomy...
$ genus
               <chr> "carni", "omni", "herbi", "omni", "herbi", "herbi", "carni", NA, "carni", "herbi", "her...
$ vore
               <chr> "Carnivora", "Primates", "Rodentia", "Soricomorpha", "Artiodactyla", "Pilosa", "Carnivo...
$ order
$ conservation <chr> "lc", NA, "nt", "lc", "domesticated", NA, "vu", NA, "domesticated", "lc", "lc", "domest...
$ sleep_total
               <dbl> 12.1, 17.0, 14.4, 14.9, 4.0, 14.4, 8.7, 7.0, 10.1, 3.0, 5.3, 9.4, 10.0, 12.5, 10.3, 8.3...
$ sleep_rem
               <dbl> NA, 1.8, 2.4, 2.3, 0.7, 2.2, 1.4, NA, 2.9, NA, 0.6, 0.8, 0.7, 1.5, 2.2, 2.0, 1.4, 3.1, ...
$ sleep_cycle
               <dbl> NA, NA, NA, 0.133, 0.667, 0.767, 0.383, NA, 0.333, NA, NA, 0.217, NA, 0.117, NA, NA, 0....
$ awake
               <dbl> 11.9, 7.0, 9.6, 9.1, 20.0, 9.6, 15.3, 17.0, 13.9, 21.0, 18.7, 14.6, 14.0, 11.5, 13.7, 1...
$ brainwt
               <dbl> NA, 0.01550, NA, 0.00029, 0.42300, NA, NA, NA, 0.07000, 0.09820, 0.11500, 0.00550, NA, ...
$ bodywt
               <dbl> 50.000, 0.480, 1.350, 0.019, 600.000, 3.850, 20.490, 0.045, 14.000, 14.800, 33.500, 0.7...
```

```
There are:
> alimpse(msleep
Observations: 83
                     11 columns
Variables: 11
$ name
               <chr> "Cheetah", "Owl monkey", "Mountain beaver", "Greater short-tailed shrew", "Cow", "Three...
               <chr> "Acinonyx", "Aotus", "Aplodontia", "Blarina", "Bos", "Bradypus", "Callorhinus", "Calomy...
$ genus
               <chr> "carni", "omni", "herbi", "omni", "herbi", "herbi", "carni", NA, "carni", "herbi", "her...
$ vore
               <chr> "Carnivora", "Primates", "Rodentia", "Soricomorpha", "Artiodactyla", "Pilosa", "Carnivo...
$ order
$ conservation <chr> "lc", NA, "nt", "lc", "domesticated", NA, "vu", NA, "domesticated", "lc", "lc", "domest...
$ sleep_total
               <dbl> 12.1, 17.0, 14.4, 14.9, 4.0, 14.4, 8.7, 7.0, 10.1, 3.0, 5.3, 9.4, 10.0, 12.5, 10.3, 8.3...
$ sleep_rem
               <dbl> NA, 1.8, 2.4, 2.3, 0.7, 2.2, 1.4, NA, 2.9, NA, 0.6, 0.8, 0.7, 1.5, 2.2, 2.0, 1.4, 3.1, ...
$ sleep_cycle
               <dbl> NA, NA, NA, 0.133, 0.667, 0.767, 0.383, NA, 0.333, NA, NA, 0.217, NA, 0.117, NA, NA, 0....
$ awake
               <dbl> 11.9, 7.0, 9.6, 9.1, 20.0, 9.6, 15.3, 17.0, 13.9, 21.0, 18.7, 14.6, 14.0, 11.5, 13.7, 1...
               <dbl> NA, 0.01550, NA, 0.00029, 0.42300, NA, NA, NA, 0.07000, 0.09820, 0.11500, 0.00550, NA, ...
$ brainwt
$ bodywt
               <dbl> 50.000, 0.480, 1.350, 0.019, 600.000, 3.850, 20.490, 0.045, 14.000, 14.800, 33.500, 0.7...
```

The names of the columns

```
The first 5 columns are
                There are:
> alimpse(msleep
                                       character variables
Observations: 83
                     11 columns
Variables: 11
$ name
               <chr> "Cheetah", "Owl monkey", "Mountain beaver", "Greater short-tailed shrew", "Cow", "Three...
               <chr> "Acinonyx", "Aotus", "Aplodontia", "Blarina", "Bos", "Bradypus", "Callorhinus", "Calomy...
$ genus
               <chr> "carni", "omni", "herbi", "omni", "herbi", "herbi", "carni", NA, "carni", "herbi", "her...
$ vore
               <chr> "Carnivora", "Primates", "Rodentia", "Soricomorpha", "Artiodactyla", "Pilosa", "Carnivo...
$ order
$ conservation <chr> "lc", NA, "nt", "lc", "domesticated", NA, "vu", NA, "domesticated", "lc", "lc", "domest...
$ sleep_total
               <dbl> 12.1, 17.0, 14.4, 14.9, 4.0, 14.4, 8.7, 7.0, 10.1, 3.0, 5.3, 9.4, 10.0, 12.5, 10.3, 8.3...
$ sleep_rem
               <dbl> NA, 1.8, 2.4, 2.3, 0.7, 2.2, 1.4, NA, 2.9, NA, 0.6, 0.8, 0.7, 1.5, 2.2, 2.0, 1.4, 3.1, ...
$ sleep_cycle
               <dbl> NA, NA, NA, 0.133, 0.667, 0.767, 0.383, NA, 0.333, NA, NA, 0.217, NA, 0.117, NA, NA, 0....
$ awake
               <dbl> 11.9, 7.0, 9.6, 9.1, 20.0, 9.6, 15.3, 17.0, 13.9, 21.0, 18.7, 14.6, 14.0, 11.5, 13.7, 1...
               <dbl> NA, 0.01550, NA, 0.00029, 0.42300, NA, NA, NA, 0.07000, 0.09820, 0.11500, 0.00550, NA, ...
$ brainwt
               <dbl> 50.000, 0.480, 1.350, 0.019, 600.000, 3.850, 20.490, 0.045, 14.000, 14.800, 33.500, 0.7...
$ bodywt
```

The names of the columns

```
The first three names of the
                                     The first 5 columns are
                There are:
> alimpse(msleep
                                                                         animals in the dataset
                     83 rows
                                       character variables
Observations: 83
                     11 columns
Variables: 11
$ name
               <chr> "Cheetah", "Owl monkey", "Mountain beaver", "Greater short-tailed shrew", "Cow", "Three...
               <chr> "Acinonyx", "Aotus", "Apiodontia", "Blarina", "Bos", "Bradypus", "Callorhinus", "Calomy...
$ genus
               <chr> "carni", "omni", "herbi", "omni", "herbi", "herbi", "carni", NA, "carni", "herbi", "her...
$ vore
               <chr> "Carnivora", "Primates", "Rodentia", "Soricomorpha", "Artiodactyla", "Pilosa", "Carnivo...
$ order
$ conservation <chr> "lc", NA, "nt", "lc", "domesticated", NA, "vu", NA, "domesticated", "lc", "lc", "domest...
$ sleep_total
               <dbl> 12.1, 17.0, 14.4, 14.9, 4.0, 14.4, 8.7, 7.0, 10.1, 3.0, 5.3, 9.4, 10.0, 12.5, 10.3, 8.3...
$ sleep_rem
               <dbl> NA, 1.8, 2.4, 2.3, 0.7, 2.2, 1.4, NA, 2.9, NA, 0.6, 0.8, 0.7, 1.5, 2.2, 2.0, 1.4, 3.1, ...
$ sleep_cycle
               <dbl> NA, NA, NA, 0.133, 0.667, 0.767, 0.383, NA, 0.333, NA, NA, 0.217, NA, 0.117, NA, NA, 0....
$ awake
               <dbl> 11.9, 7.0, 9.6, 9.1, 20.0, 9.6, 15.3, 17.0, 13.9, 21.0, 18.7, 14.6, 14.0, 11.5, 13.7, 1...
               <dbl> NA, 0.01550, NA, 0.00029, 0.42300, NA, NA, NA, 0.07000, 0.09820, 0.11500, 0.00550, NA, ...
$ brainwt
               <dbl> 50.000, 0.480, 1.350, 0.019, 600.000, 3.850, 20.490, 0.045, 14.000, 14.800, 33.500, 0.7...
$ bodywt
```

The names of the columns

Equivalent to:

```
msleep %>%
                                            filter(msleep, order == "Primates")
    filter(order == "Primates")
# A tibble: 12 x 11
                                      order
                                             conservation sleep_total sleep_rem sleep_cycle awake
                                                                                                     brainwt bodvwt
   name
                   genus
                                vore
  <chr>
                   <chr>
                                <chr> <chr>
                                             <chr>
                                                                 <dbl>
                                                                           <dbl>
                                                                                        <dbl> <dbl>
                                                                                                        < dbl>
                                                                                                               <dbl>
                                                                           1.80
                                                                                               7.00
1 Owl monkey
                   Aotus
                                omni
                                      Prima... <NA>
                                                                 17.0
                                                                                       NA
                                                                                                     0.0155
                                                                                                               0.480
2 Grivet
                   Cercopithe... omni
                                      Prima...
                                                                 10.0
                                                                           0.700
                                                                                       NA
                                                                                              14.0
                                                                                                    NA
                                                                                                               4.75
                                             lc
 3 Patas monkey
                   Erythroceb... omni
                                      Prima...
                                                                 10.9
                                                                           1.10
                                                                                       NA
                                                                                              13.1
                                                                                                     0.115
                                                                                                              10.0
                                             lc
4 Galago
                   Galago
                                      Prima... <NA>
                                                                  9.80
                                                                           1.10
                                                                                        0.550 14.2
                                                                                                     0.00500
                                                                                                               0.200
                                omni
 5 Human
                                      Prima... <NA>
                                                                  8.00
                                                                            1.90
                                                                                        1.50
                                                                                              16.0
                                                                                                     1.32
                                                                                                              62.0
                   Homo
                                omni
                                herbi Prima... vu
 6 Mongoose lemur
                   Lemur
                                                                  9.50
                                                                           0.900
                                                                                       NA
                                                                                              14.5
                                                                                                    NA
                                                                                                               1.67
                                      Prima... <NA>
                                                                 10.1
                                                                           1.20
                                                                                        0.750 13.9
                                                                                                     0.179
                                                                                                               6.80
7 Macaque
                   Macaca
                                omni
8 Slow loris
                   Nyctibeus
                                carni Prima... <NA>
                                                                 11.0
                                                                                              13.0
                                                                                                     0.0125
                                                                                                               1.40
                                                                          NA
                                                                                       NA
9 Chimpanzee
                                      Prima... <NA>
                                                                  9.70
                                                                           1.40
                                                                                        1.42 14.3
                                                                                                     0.440
                                                                                                              52.2
                   Pan
                                omni
10 Baboon
                                      Prima... <NA>
                                                                           1.00
                                                                                        0.667 14.6
                                                                                                     0.180
                                                                                                              25.2
                   Papio
                                omni
                                                                  9.40
11 Potto
                   Perodictic... omni
                                      Prima...lc
                                                                 11.0
                                                                          NA
                                                                                       NA
                                                                                              13.0
                                                                                                    NA
                                                                                                               1.10
12 Squirrel monkey Saimiri
                                                                  9.60
                                                                                              14.4
                                                                                                     0.0200
                                                                                                               0.743
                                omni
                                      Prima... <NA>
                                                                           1.40
                                                                                       NA
```

```
> msleep %>%
   filter(order == "Primates", sleep_total > 10)
# A tibble: 5 x 11
                                           conservation sleep_total sleep_rem sleep_cycle awake brainwt bodywt
                                 order
  name
               genus
                            vore
                            <chr> <chr>
                                           <chr>
                                                              <dbl>
                                                                        <dbl>
                                                                                    <dbl> <dbl>
                                                                                                  <dbl> <dbl>
  <chr>>
               <chr>
1 Owl monkey
                                 Primates <NA>
                                                                                                         0.480
               Aotus
                                                               17.0
                                                                         1.80
                                                                                            7.00
                                                                                                 0.0155
                            omni
                                                                                   NA
                                 Primates lc
                                                                                                 0.115
2 Patas monkey Erythrocebus omni
                                                               10.9
                                                                         1.10
                                                                                           13.1
                                                                                                        10.0
                                                                                   NA
               Macaca
                                 Primates <NA>
                                                               10.1
                                                                         1.20
                                                                                    0.750 13.9
                                                                                                  0.179
3 Macaque
                            omni
                                                                                                          6.80
4 Slow loris Nyctibeus
                            carni Primates <NA>
                                                               11.0
                                                                        NA
                                                                                   NA
                                                                                           13.0
                                                                                                 0.0125
                                                                                                         1.40
               Perodicticus omni
5 Potto
                                  Primates lc
                                                               11.0
                                                                        NA
                                                                                   NA
                                                                                           13.0
                                                                                                NA
                                                                                                          1.10
```

```
Gives the same results: msleep %>%
filter(order == "Primates" & sleep_total > 10)
```

```
> msleep %>%
   filter(order == "Primates", sleep_total > 10) %>%
   select(name, sleep_total, sleep_rem, sleep_cycle)
# A tibble: 5 x 4
              sleep_total sleep_rem sleep_cycle
 name
                                       <dbl>
 <chr>
                   <dbl>
                            <dbl>
                    17.0
                             1.80
1 Owl monkey
                                       NA
                    10.9
                             1.10
2 Patas monkey
                                       NA
                    10.1 1.20
3 Macaque
                                      0.750
4 Slow loris
                    11.0
                            NA
                                       NA
5 Potto
                    11.0
                            NA
                                       NA
```

```
> msleep %>%
   filter(order == "Primates", sleep_total > 10) %>%
   select(name, total=sleep_total, rem=sleep_rem, cycle=sleep_cycle)
# A tibble: 5 x 4
              total
                          cycle
 name
                     rem
             <dbl> <dbl> <dbl>
 <chr>
1 Owl monkey 17.0 1.80 NA
2 Patas monkey 10.9 1.10 NA
          10.1 1.20 0.750
3 Macaque
4 Slow loris 11.0 NA
                         NA
5 Potto 11.0 NA
                      NA
```

```
> msleep %>%
      filter(order == "Primates", sleep_total > 10) %>%
     rename(total=sleep_total, rem=sleep_rem, cycle=sleep_cycle)
# A tipple: 5 x 11
                                           conservation total
                                                                rem cycle awake brainwt bodywt
                            vore order
  name
               aenus
  <chr>
               <chr>
                            <chr> <chr>
                                           <chr>
                                                        <dbl> <dbl> <dbl> <dbl>
                                                                                  <dbl> <dbl>
1 Owl monkey
                                 Primates <NA>
                                                         17
                                                                1.8 NA
                                                                                 0.0155
                                                                                          0.48
               Aotus
                            omni
                                                                            7
2 Patas monkey Erythrocebus omni
                                 Primates lc
                                                         10.9
                                                                1.1 NA
                                                                           13.1
                                                                                 0.115
                                                                                         10
                                                         10.1
                                                                1.2 0.75
                                                                           13.9
                                                                                          6.8
3 Macaque
               Macaca
                            omni
                                  Primates <NA>
                                                                                 0.179
4 Slow loris
              Nyctibeus
                            carni Primates <NA>
                                                         11
                                                               NA
                                                                    NA
                                                                           13
                                                                                 0.0125
                                                                                          1.4
5 Potto
               Perodicticus omni Primates lc
                                                         11
                                                                           13
                                                               NA
                                                                    NA
                                                                                NA
                                                                                          1.1
```

Reordering Data

```
> msleep %>%
   filter(order == "Primates", sleep_total > 10) %>%
   select(name, sleep_rem, sleep_cycle, sleep_total)
# A tibble: 5 x 4
              sleep_rem sleep_cycle sleep_total
 name
 <chr>
                 <dbl>
                             <dbl>
                                        <dbl>
1 Owl monkey
                  1.80
                                         17.0
                            NA
               1.10
2 Patas monkey
                            NA
                                         10.9
              1.20
3 Macaque
                           0.750
                                         10.1
4 Slow loris
                 NA
                            NA
                                         11.0
                                         11.0
5 Potto
                 NA
                            NA
```

```
> msleep %>%
    filter(order == "Primates", sleep_total > 10) %>%
    select(name, sleep_rem, sleep_cycle, sleep_total) %>%
   arrange(sleep_total)
# A tibble: 5 x 4
               sleep_rem sleep_cycle sleep_total
  name
                   <dbl>
                               <dbl>
                                            <dbl>
  <chr>
                                                   smallest
                                             10.1
1 Macague
                    1.20
                               0.750
                    1.10
                                             10.9
2 Patas monkey
                               NA
3 Slow loris
                               NA
                                             11.0
                   NA
                               NA
                                             11.0
4 Potto
                   NA
                    1.80
                               NA
                                             17.0
5 Owl monkey
                                                   largest
```

```
> msleep %>%
   filter(order == "Primates", sleep_total > 10) %>%
   select(name, sleep_rem, sleep_cycle, sleep_total) %>%
   arrange(desc(sleep_total))
# A tibble: 5 x 4
               sleep_rem sleep_cycle sleep_total
  name
                   <dbl>
                              <dbl>
                                           <dbl>
  <chr>
                                                 largest
                                           17.0
1 Owl monkey
                   1.80
                             NA
2 Slow loris
                  NA
                              NA
                                            11.0
3 Potto
                  NA
                              NA
                                           11.0
                1.10
                                           10.9
4 Patas monkey
                              NA
5 Macaque
                    1.20
                              0.750
                                            10.1
```

```
> msleep %>%
   filter(order == "Primates", sleep_total > 10) %>%
   select(name, sleep_rem, sleep_cycle, sleep_total) %>%
   arrange(name)
# A tibble: 5 x 4
               sleep_rem sleep_cycle sleep_total
  name
                  <dbl>
                              <dbl>
                                          <dbl>
 <chr>
1 Macaque
                   1.20
                              0.750
                                           10.1
2 Owl monkey
                   1.80
                             NA
                                           17.0
3 Patas monkey
                   1.10
                             NA
                                           10.9
4 Potto
                             NA
                                           11.0
                  NA
5 Slow loris
                             NA
                                            11.0
                  NA
```

Sorted alphabetically by name

Sort alphabetically by name, then total sleep:

arrange(name, sleep_total)

Manipulating Data

```
> msleep %>%
   filter(order == "Primates", sleep_total > 10) %>%
    select(name, sleep_rem, sleep_cycle, sleep_total) %>%
    arrange(name) %>%
   mutate(sleep_total_min = sleep_total * 60)
# A tibble: 5 x 5
               sleep_rem sleep_cycle sleep_total sleep_total_min
  name
                                            <dbl>
  <chr>
                   <dbl>
                               <dbl>
                                                            <dbl>
                               0.750
1 Macaque
                    1.20
                                             10.1
                                                              606
2 Owl monkey
                                             17.0
                                                             1020
                    1.80
                              NA
3 Patas monkey
                1.10
                              NA
                                             10.9
                                                              654
                                             11.0
                                                              660
4 Potto
                   NA
                              NA
5 Slow loris
                                             11.0
                   NA
                              NA
                                                              660
```

A whole new column!

```
## if not already installed, you'll have to run the following line of code
install.packages('httr')

## load the library
library("httr")

## download file
GET("https://raw.githubusercontent.com/suzanbaert/RTutorials/master/Rmd_originals/conservation_explanation.csv",
write_disk(tf <- tempfile(fileext = ".csv")))
conservation <- read_csv(tf)

## take a look at this file
head(conservation)</pre>
```

```
> conservation
# A tibble: 11 x 1
                                     Tidy data violation!
   `conservation abbreviation`
                                     Space in column name should
                                     be an underscore.
   <chr>
 1 EX = Extinct
                                   Tidy data violation!
 2 EW = Extinct in the wild
                                   There are two pieces of
 3 CR = Critically Endangered
                                   information in a column: (1)
4 EN = Endangered
                                   abbreviation and (2) description.
 5 VU = Vulnerable
 6 NT = Near Threatened
 7 LC = Least Concern
 8 DD = Data deficient
 9 NE = Not evaluated
10 PE = Probably extinct (informal)
11 PEW = Probably extinct in the wild (informal)
```

```
> conservation %>%
    separate(`conservation abbreviation`,
             into = c("abbreviation", "description"), sep = " = ")
# A tibble: 11 x 2
   abbreviation description
                <chr>
   <chr>
                Extinct
 1 EX
 2 EW
                Extinct in the wild
 3 CR
                Critically Endangered
 4 EN
                Endangered
 5 VU
                Vulnerable
 6 NT
                Near Threatened
 7 LC
                Least Concern
 8 DD
                Data deficient
                Not evaluated
 9 NE
10 PE
                Probably extinct (informal)
11 PEW
                Probably extinct in the wild (informal)
```

```
> conservation %>%
    separate(`conservation abbreviation`,
             into = c("abbreviation", "description"), sep = " = ") %>%
    unite(united_col, abbreviation, description, sep=" = ")
 A tibble: 11 x 1
   united_col
  <chr>
 1 EX = Extinct
 2 EW = Extinct in the wild
 3 CR = Critically Endangered
 4 EN = Endangered
 5 VU = Vulnerable
 6 NT = Near Threatened
 7 LC = Least Concern
 8 DD = Data deficient
 9 NE = Not evaluated
10 PE = Probably extinct (informal)
11 PEW = Probably extinct in the wild (informal)
```

```
> conservation %>%
   clean_names()
# A tibble: 11 x 1
   conservation_abbreviation
   <chr>>
                    Adds underscore to column name
1 FX = Fxtinct
2 EW = Extinct in the wild
3 CR = Critically Endangered
4 EN = Endangered
 5 VU = Vulnerable
 6 NT = Near Threatened
 7 LC = Least Concern
 8 DD = Data deficient
9 NE = Not evaluated
10 PE = Probably extinct (informal)
11 PEW = Probably extinct in the wild (informal)
```

```
> msleep %>%
    mutate(conservation = toupper(conservation)) %>%
   left_join(conserve, by = c("conservation" = "abbreviation"))
# A tibble: 83 x 12
             genus vore order conservation sleep_total sleep_rem sleep_cycle awake brainwt bodywt description
   name
   <chr> <chr> <chr> <chr> <chr> <chr>
                                                        <dbl>
                                                                   <dbl>
                                                                                <dbl> <dbl>
                                                                                                 <dbl> <dbl> <chr>
                                                        12.1
                                                                  NA
                                                                                       11.9 NA
                                                                                                        5.00e<sup>+1</sup> Least Conc...
 1 Cheetah Acino... carni Carn... LC
                                                                               NA
 2 Owl mon... Aotus omni Prim... <NA>
                                                        17.0
                                                                                        7.00 1.55e^{-2} 4.80e^{-1} <NA>
                                                                   1.80
                                                                               NA
 3 Mountai... Aplod... herbi Rode... NT
                                                        14.4
                                                                   2.40
                                                                                        9.60 NA
                                                                                                        1.35e<sup>+0</sup> Near Threa...
                                                                               NA
                                                                                0.133 9.10 2.90e<sup>-4</sup> 1.90e<sup>-2</sup> Least Conc...
 4 Greater... Blari... omni Sori... LC
                                                        14.9
                                                                   2.30
                     herbi Arti... DOMESTICATED
                                                                   0.700
                                                                                0.667 20.0
                                                                                               4.23e^{-1} 6.00e^{+2} <NA>
 5 Cow
             Bos
                                                         4.00
 6 Three-t... Brady... herbi Pilo... <NA>
                                                                                                        3.85e^{+0} <NA>
                                                        14.4
                                                                   2.20
                                                                                0.767 9.60 NA
 7 Norther... Callo... carni Carn... VU
                                                         8.70
                                                                   1.40
                                                                                0.383 15.3 NA
                                                                                                        2.05e<sup>+1</sup> Vulnerable
 8 Vesper ... Calom... <NA> Rode... <NA>
                                                                                       17.0 NA
                                                                                                        4.50e^{-2} <NA>
                                                         7.00
                                                                  NA
                                                                               NA
 9 Dog
             Canis carni Carn... DOMESTICATED
                                                        10.1
                                                                   2.90
                                                                                0.333 \ 13.9 \ 7.00e^{-2} \ 1.40e^{+1} < NA>
10 Roe deer Capre... herbi Arti... LC
                                                                  NA
                                                                                       21.0
                                                                                             9.82e<sup>-2</sup> 1.48e<sup>+1</sup> Least Conc...
                                                         3.00
                                                                               NA
# ... with 73 more rows
```

Summarizing Data

```
> msleep
# A tibble: 83 x 11
   name
                     genus
                              vore order
                                               conservation sleep_total sleep_rem sleep_cycle awake brainwt bodywt
   <chr>
                     <chr>
                               <chr> <chr>
                                               <chr>
                                                                    <dbl>
                                                                               <dbl>
                                                                                             <dbl> <dbl>
                                                                                                             <dbl> <dbl>
 1 Cheetah
                     Acinonyx carni Carnivo... lc
                                                                    12.1
                                                                              NA
                                                                                                   11.9 NA
                                                                                                                    5.00e+1
 2 Owl monkey
                     Aotus
                               omni Primates <NA>
                                                                    17.0
                                                                               1.80
                                                                                                    7.00 1.55e<sup>-2</sup> 4.80e<sup>-1</sup>
 3 Mountain beaver Aplodon... herbi Rodentia nt
                                                                    14.4
                                                                               2.40
                                                                                                    9.60 NA
                                                                                                                   1.35e+0
 4 Greater short-... Blarina omni Soricom... lc
                                                                    14.9
                                                                               2.30
                                                                                            0.133 9.10 2.90e<sup>-4</sup> 1.90e<sup>-2</sup>
 5 Cow
                     Bos
                              herbi Artioda... domesticated
                                                                     4.00
                                                                               0.700
                                                                                            0.667 20.0
                                                                                                          4.23e<sup>-1</sup> 6.00e<sup>+2</sup>
                                                                                                                    3.85e+0
 6 Three-toed slo... Bradypus herbi Pilosa <NA>
                                                                    14.4
                                                                               2.20
                                                                                            0.767 9.60 NA
                                                                                            0.383 15.3 NA
 7 Northern fur s... Callorh... carni Carnivo... vu
                                                                     8.70
                                                                               1.40
                                                                                                                   2.05e+1
 8 Vesper mouse
                    Calomys <NA> Rodentia <NA>
                                                                     7.00
                                                                              NA
                                                                                                   17.0 NA
                                                                                                                    4.50e<sup>-2</sup>
 9 Dog
                    Canis
                              carni Carnivo... domesticated
                                                                    10.1
                                                                               2.90
                                                                                            0.333 13.9 7.00e<sup>-2</sup> 1.40e<sup>+1</sup>
10 Roe deer
                    Capreol... herbi Artioda... lc
                                                                     3.00
                                                                                                   21.0
                                                                                                           9.82e<sup>-2</sup> 1.48e<sup>+1</sup>
                                                                              NA
# ... with 73 more rows
> msleep %>%
  group_by(order)
# A tibble: 83 x 11
# Groups: order [19]
```

	name	genus	vore	order	conservation	sleep_total	sleep_rem	sleep_cycle	awake	brainwt	bodywt
	<chr></chr>	<chr></chr>	<chr></chr>	<chr></chr>	<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
1	Cheetah	Acinonyx	carni	Carnivo	lc	12.1	NA	NA	11.9	NA	5.00e ⁺¹
2	Owl monkey	Aotus	omni	Primates	<na></na>	17.0	1.80	NA	7.00	1.55e ⁻²	4.80e ⁻¹
3	Mountain beaver	Aplodon	herbi	Rodentia	nt	14.4	2.40	NA	9.60	NA	1.35e ⁺⁰
4	Greater short	Blarina	omni	Soricom	lc	14.9	2.30	0.133	9.10	2.90e ⁻⁴	1.90e ⁻²
5	Cow	Bos	herbi	Artioda	domesticated	4.00	0.700	0.667	20.0	4.23e ⁻¹	6.00e ⁺²
6	Three-toed slo	Bradypus	herbi	Pilosa	<na></na>	14.4	2.20	0.767	9.60	NA	3.85e ⁺⁰
7	Northern fur s	Callorh	carni	Carnivo	vu	8.70	1.40	0.383	15.3	NA	2.05e ⁺¹
8	Vesper mouse	Calomys	<na></na>	Rodentia	<na></na>	7.00	NA	NA	17.0	NA	4.50e ⁻²
9	Dog	Canis	carni	Carnivo	domesticated	10.1	2.90	0.333	13.9	7.00e ⁻²	1.40e ⁺¹
10	Roe deer	Capreol	herbi	Artioda	lc	3.00	NA	NA	21.0	9.82e ⁻²	1.48e ⁺¹
#	with 73 more	NOWE									

... with /3 more rows

```
> msleep %>%
+ select(order) %>%
+ summarize(N=n())
# A tibble: 1 x 1
          N
          <int>
1 83
          Same as nrow (msleep)
```

```
> msleep %>%
   group_by(order) %>%
    select(order) %>%
   summarize(N=n())
# A tibble: 19 x 2
   order
                       N
   <chr>>
                   <int>
 1 Afrosoricida
 2 Artiodactyla
                       6
 3 Carnivora
                      12
 4 Cetacea
                       3
 5 Chiroptera
 6 Cingulata
                       2
 7 Didelphimorphia
 8 Diprotodontia
                       2
9 Erinaceomorpha
10 Hyracoidea
11 Lagomorpha
12 Monotremata
13 Perissodactyla
                       3
14 Pilosa
15 Primates
                      12
16 Proboscidea
17 Rodentia
                      22
18 Scandentia
                       1
19 Soricomorpha
                       5
```

```
> msleep %>%
    group_by(order) %>%
    select(order, sleep_total) %>%
   summarize(N=n(), mean_sleep=mean(sleep_total))
# A tibble: 19 x 3
   order
                       N mean_sleep
   <chr>>
                   <int>
                               <dbl>
 1 Afrosoricida
                               15.6
                                4.52
 2 Artiodactyla
 3 Carnivora
                      12
                               10.1
                                4.50
 4 Cetacea
 5 Chiroptera
                               19.8
 6 Cingulata
                               17.8
 7 Didelphimorphia
                               18.7
 8 Diprotodontia
                               12.4
 9 Erinaceomorpha
                               10.2
10 Hyracoidea
                                5.67
11 Lagomorpha
                                8.40
12 Monotremata
                                8.60
                                3.47
13 Perissodactyla
14 Pilosa
                               14.4
                      12
15 Primates
                               10.5
16 Proboscidea
                                3.60
                       22
17 Rodentia
                               12.5
18 Scandentia
                                8.90
19 Soricomorpha
                               11.1
```

```
> msleep %>%
    tabyl(order)
             order
                    n percent
      Afrosoricida
                       0.0120
      Artiodactyla
                       0.0723
         Carnivora 12
                       0.1446
                       0.0361
           Cetacea
5
        Chiroptera
                       0.0241
6
         Cingulata
                       0.0241
   Didelphimorphia
                       0.0241
8
     Diprotodontia
                        0.0241
9
    Erinaceomorpha
                       0.0241
10
        Hyracoidea
                       0.0361
11
        Lagomorpha
                       0.0120
12
       Monotremata
                        0.0120
13
    Perissodactyla
                       0.0361
            Pilosa
14
                       0.0120
15
          Primates 12
                       0.1446
16
       Proboscidea
                       0.0241
17
          Rodentia 22
                       0.2651
18
        Scandentia
                       0.0120
19
      Soricomorpha
                       0.0602
```

> summary(msleep\$awake)

Min. 1st Qu. Median

4.1 10.2 13.9

Mean 3rd Qu. Max.

13.6 16.1 22.1

> skim(msleep)

Skim summary statistics

n obs: 83

n variables: 11

Variable type: character

variable	missing	complete	n	min	max	empty	n_unique
conservation	29	54	83	2	12	0	6
genus	0	83	83	3	13	0	77
name	0	83	83	3	30	0	83
order	0	83	83	6	15	0	19
vore	7	76	83	4	7	0	4

Variable type: numeric

variable	missing	complete	n	mean	sd	p0	p25	p50	p75	p100	hist
awake	0	83	83	13.57	4.45	4.1	10.25	13.9	16.15	22.1	
bodywt	0	83	83	166.14	786.84	0.005	0.17	1.67	41.75	6654	
brainwt	27	56	83	0.28	0.98	0.00014	0.0029	0.012	0.13	5.71	
sleep_cycle	51	32	83	0.44	0.36	0.12	0.18	0.33	0.58	1.5	
sleep_rem	22	61	83	1.88	1.3	0.1	0.9	1.5	2.4	6.6	
sleep_total	0	83	83	10.43	4.45	1.9	7.85	10.1	13.75	19.9	

Filtering, Re-ordering, Manipulating, and Summarizing