

Section 2. Exploring Data Frame

Sandeep Kumar

March 16, 2018

Lets explore mtcars dataset available in R package.

```
class(mtcars)
```

```
## [1] "data.frame"
```

Dimensions

```
dim(mtcars)
```

```
## [1] 32 11
```

Names

```
names(mtcars)
```

```
## [1] "mpg" "cyl" "disp" "hp" "drat" "wt" "qsec" "vs" "am" "gear"  
## [11] "carb"
```

Structure

```
str(mtcars)
```

```
## 'data.frame': 32 obs. of 11 variables:  
## $ mpg : num 21 21 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 ...  
## $ cyl : num 6 6 4 6 8 6 8 4 4 6 ...  
## $ disp: num 160 160 108 258 360 ...  
## $ hp : num 110 110 93 110 175 105 245 62 95 123 ...  
## $ drat: num 3.9 3.9 3.85 3.08 3.15 2.76 3.21 3.69 3.92 3.92 ...  
## $ wt : num 2.62 2.88 2.32 3.21 3.44 ...  
## $ qsec: num 16.5 17 18.6 19.4 17 ...  
## $ vs : num 0 0 1 1 0 1 0 1 1 1 ...  
## $ am : num 1 1 1 0 0 0 0 0 0 0 ...  
## $ gear: num 4 4 4 3 3 3 3 4 4 4 ...  
## $ carb: num 4 4 1 1 2 1 4 2 2 4 ...
```

```
head(mtcars)
```

```
##           mpg cyl disp  hp drat   wt  qsec vs am gear carb  
## Mazda RX4      21.0   6  160 110 3.90 2.620 16.46  0  1    4    4  
## Mazda RX4 Wag  21.0   6  160 110 3.90 2.875 17.02  0  1    4    4  
## Datsun 710     22.8   4  108  93 3.85 2.320 18.61  1  1    4    1  
## Hornet 4 Drive 21.4   6  258 110 3.08 3.215 19.44  1  0    3    1
```

```
## Hornet Sportabout 18.7    8   360 175 3.15 3.440 17.02  0  0    3    2
## Valiant            18.1    6   225 105 2.76 3.460 20.22  1  0    3    1
```

```
tail(mtcars)
```

```
##           mpg  cyl  disp  hp drat    wt  qsec vs am gear carb
## Porsche 914-2 26.0   4 120.3  91 4.43  2.140 16.7  0  1    5    2
## Lotus Europa  30.4   4  95.1 113 3.77  1.513 16.9  1  1    5    2
## Ford Pantera L 15.8   8 351.0 264 4.22  3.170 14.5  0  1    5    4
## Ferrari Dino   19.7   6 145.0 175 3.62  2.770 15.5  0  1    5    6
## Maserati Bora   15.0   8 301.0 335 3.54  3.570 14.6  0  1    5    8
## Volvo 142E     21.4   4 121.0 109 4.11  2.780 18.6  1  1    4    2
```

```
summary(mtcars)
```

```
##           mpg           cyl           disp           hp
##  Min.   :10.40   Min.   :4.000   Min.   : 71.1   Min.   : 52.0
## 1st Qu.:15.43   1st Qu.:4.000   1st Qu.:120.8   1st Qu.: 96.5
##  Median :19.20   Median :6.000   Median :196.3   Median :123.0
##   Mean   :20.09   Mean   :6.188   Mean   :230.7   Mean   :146.7
## 3rd Qu.:22.80   3rd Qu.:8.000   3rd Qu.:326.0   3rd Qu.:180.0
##   Max.   :33.90   Max.   :8.000   Max.   :472.0   Max.   :335.0
##           drat           wt           qsec           vs
##  Min.   :2.760   Min.   :1.513   Min.   :14.50   Min.   :0.0000
## 1st Qu.:3.080   1st Qu.:2.581   1st Qu.:16.89   1st Qu.:0.0000
##  Median :3.695   Median :3.325   Median :17.71   Median :0.0000
##   Mean   :3.597   Mean   :3.217   Mean   :17.85   Mean   :0.4375
## 3rd Qu.:3.920   3rd Qu.:3.610   3rd Qu.:18.90   3rd Qu.:1.0000
##   Max.   :4.930   Max.   :5.424   Max.   :22.90   Max.   :1.0000
##           am           gear           carb
##  Min.   :0.0000   Min.   :3.000   Min.   :1.000
## 1st Qu.:0.0000   1st Qu.:3.000   1st Qu.:2.000
##  Median :0.0000   Median :4.000   Median :2.000
##   Mean   :0.4062   Mean   :3.688   Mean   :2.812
## 3rd Qu.:1.0000   3rd Qu.:4.000   3rd Qu.:4.000
##   Max.   :1.0000   Max.   :5.000   Max.   :8.000
```

The above data frame just had numeric data. Lets understand another dataset

```
head(chickwts)
```

```
##   weight      feed
## 1    179 horsebean
## 2    160 horsebean
## 3    136 horsebean
## 4    227 horsebean
## 5    217 horsebean
## 6    168 horsebean
```

```
str(chickwts)
```

```
## 'data.frame':    71 obs. of  2 variables:
## $ weight: num  179 160 136 227 217 168 108 124 143 140 ...
## $ feed   : Factor w/ 6 levels "casein","horsebean",...: 2 2 2 2 2 2 2 2 2 2
## ...
```

Data Frames in other packages - eg mpg or diamonds in ggplot2

```
#install.packages("ggplot2")
```

```
library(ggplot2)
```

```
## Warning: package 'ggplot2' was built under R version 3.4.3
```

```
head(diamonds)
```

```
## # A tibble: 6 x 10
```

```
##   carat      cut color clarity depth table price      x      y      z
##   <dbl>     <ord> <ord>   <ord> <dbl> <dbl> <int> <dbl> <dbl> <dbl>
## 1  0.23     Ideal   E     SI2  61.5   55  326  3.95  3.98  2.43
## 2  0.21  Premium   E     SI1  59.8   61  326  3.89  3.84  2.31
## 3  0.23     Good    E     VS1  56.9   65  327  4.05  4.07  2.31
## 4  0.29  Premium   I     VS2  62.4   58  334  4.20  4.23  2.63
## 5  0.31     Good    J     SI2  63.3   58  335  4.34  4.35  2.75
## 6  0.24 Very Good   J     VVS2  62.8   57  336  3.94  3.96  2.48
```

```
str(diamonds)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':    53940 obs. of  10 variables:
## $ carat : num  0.23 0.21 0.23 0.29 0.31 0.24 0.24 0.26 0.22 0.23 ...
## $ cut : Ord.factor w/ 5 levels "Fair"<"Good"<...: 5 4 2 4 2 3 3 3 1 3 .
## ..
## $ color : Ord.factor w/ 7 levels "D"<"E"<"F"<"G"<...: 2 2 2 6 7 7 6 5 2 5
## ...
## $ clarity: Ord.factor w/ 8 levels "I1"<"SI2"<"SI1"<...: 2 3 5 4 2 6 7 3 4
## 5 ...
## $ depth : num  61.5 59.8 56.9 62.4 63.3 62.8 62.3 61.9 65.1 59.4 ...
## $ table : num  55 61 65 58 58 57 57 55 61 61 ...
## $ price : int  326 326 327 334 335 336 336 337 337 338 ...
## $ x : num  3.95 3.89 4.05 4.2 4.34 3.94 3.95 4.07 3.87 4 ...
## $ y : num  3.98 3.84 4.07 4.23 4.35 3.96 3.98 4.11 3.78 4.05 ...
## $ z : num  2.43 2.31 2.31 2.63 2.75 2.48 2.47 2.53 2.49 2.39 ...
```

```
ggplot2::mpg
```

```
## # A tibble: 234 x 11
```

```
##   manufacturer      model displ  year   cyl    trans  drv   cty   hwy
##   <chr>          <chr> <dbl> <int> <int>   <chr> <chr> <int> <int>
## 1      audi         a4    1.8  1999     4 auto(l5) f     18    29
## 2      audi         a4    1.8  1999     4 manual(m5) f     21    29
## 3      audi         a4    2.0  2008     4 manual(m6) f     20    31
## 4      audi         a4    2.0  2008     4 auto(av) f     21    30
## 5      audi         a4    2.8  1999     6 auto(l5) f     16    26
## 6      audi         a4    2.8  1999     6 manual(m5) f     18    26
```

```
## 7      audi      a4  3.1 2008      6  auto(av)      f    18    27
## 8      audi a4 quattro  1.8 1999      4 manual(m5)    4    18    26
## 9      audi a4 quattro  1.8 1999      4  auto(l5)     4    16    25
## 10     audi a4 quattro  2.0 2008      4 manual(m6)    4    20    28
## # ... with 224 more rows, and 2 more variables: fl <chr>, class <chr>
```

```
data(package = "ggplot2")
```