

5.1 Data Introspection

Contents

Motivation	1
Read the raw data again	1
Printing the head of the dataframe	1
Install tibble	4
dimensions of the data frame: dim	5
names of the columns: names	5
summary of the data: summary	6
structure of the data: str	8
data types: typeof	9
using sapply , length , sort	9
An inventory of the kind of data we have: table	10

Motivation

This section is about getting familiar with our data. We will be using functions to know the size of our table or data frame, the names of the columns or variables, the structure of the data and the type of data for each of the variables or columns.

Read the raw data again

```
# load the library xlsx
library(xlsx)

# read the raw data
myXl <- read.xlsx("../extdata/oilfield_100w_raw_data.xlsx",
                  sheetIndex = 1, stringsAsFactors = FALSE)
```

Printing the head of the dataframe

Let's print 6 rows of data with the function `head()`. You will see a long printing. We will fix this in a minute. Read on.

```
# the function head() prints the first 6 rows
# to print the last 6 rows use tail()
print(head(myXl))
```

	Wellname	Company	Analyst	Field	Location	Platform	Fluid
#:> 1	PSCO-M005-TS	Oil Gains Co.	Aida	PISCO	M005-TS	M	0
#:> 2	PSCO-M007-TS	Oil Gains Co.	Aida	PISCO	M007-TS	M	0
#:> 3	PSCO-M004-Ls	Oil Gains Co.	Aida	PISCO	M004-Ls	M	0
#:> 4	PSCO-M008-TS	Oil Gains Co.	Aida	PISCO	M008-TS	M	0
#:> 5	PSCO-M010-SS	Oil Gains Co.	Aida	PISCO	M010-SS	M	0
#:> 6	PSCO-M006-TS	Oil Gains Co.	Aida	PISCO	M006-TS	M	0

```
#:> WellType AL_Method Completion SandControl WT_COUNT PVT_GOR PVT_API
#:> 1      0      1      0      0      27 445.7 36.0
#:> 2      0      1      0      0      22 473.0 36.0
```

[illegible]

[illegible]

[illegible]

It looks pretty long. Let's try with a package that adds better printing capabilities: `tibble`.

Install tibble

Install it with `install.packages("tibble")`

```
library(tibble)           # load the package
myXl <- as_tibble(myXl)   # convert the data frame to a tibble

head(myXl)

#> # A tibble: 6 x 51
#>   Wellname      Company Analyst Field Location Platform Fluid
#>   <chr>        <chr>    <chr> <chr>    <chr>    <chr> <dbl>
#> 1 PSCO-M005-TS Oil Gains Co.  Aida PISCO M005-TS      M      0
#> 2 PSCO-M0007-TS Oil Gains Co.  Aida PISCO M007-TS      M      0
#> 3 PSCO-M004-LS Oil Gains Co.  Aida PISCO M004-LS      M      0
#> 4 PSCO-M008-TS Oil Gains Co.  Aida PISCO M008-TS      M      0
#> 5 PSCO-M010-SS Oil Gains Co.  Aida PISCO M010-SS      M      0
#> 6 PSCO-M006-TS Oil Gains Co.  Aida PISCO M006-TS      M      0
#> # ... with 44 more variables: WellType <dbl>, AL_Method <dbl>,
#> # Completion <dbl>, SandControl <dbl>, WT_COUNT <dbl>, PVT_GOR <dbl>,
#> # PVT_API <dbl>, PVT_SG_gas <dbl>, PVT_WaterSalinity <dbl>,
#> # PVT_H2S <dbl>, PVT_CO2 <dbl>, PVT_PB_CORR <dbl>, PVT_VISC_CORR <dbl>,
#> # PVT_BPTTEMP <dbl>, PVT_BPPRES <dbl>, VLP_CORR <dbl>, IPR_CORR <dbl>,
#> # IPR_RESPRES <dbl>, IPR_RESTEMP <dbl>, IPR_TOTGOR <dbl>, IPR_WC <dbl>,
#> # IPR_VOGELRT <dbl>, IPR_VOGELPRES <dbl>, IPR_PI <dbl>, GEO_THMD <chr>,
#> # GEO_THTEMP <chr>, GL_method <dbl>, GL_ArrayMandrels <chr>,
#> # GL Vdepth <dbl>, GL GSG <dbl>, GL CO2 <dbl>, WT DATE <chr>.
```

```
#:> # WT_THT <chr>, WT_LIQRT <chr>, WT_WC <chr>, WT_THP <chr>, WT_GOR <chr>,
#:> # WT_GLIR <chr>, WT_DEPTH <chr>, WT_Enable <chr>, WT_GDEPTH <chr>,
#:> # WT_GPRES <chr>, WT_RESPRES <chr>, ProsperFilename <chr>
```

```
tail(myXl)
#:> # A tibble: 6 x 51
#:>       Wellname      Company Analyst Field Location Platform Fluid
#:>       <chr>         <chr>   <chr> <chr>   <chr>   <chr> <dbl>
#:> 1 PSCO-S021-TS Oil Gains Co. Camden PISCO S021-TS      S      0
#:> 2 PSCO-S016-SS Oil Gains Co. Camden PISCO S016-SS      S      0
#:> 3 PSCO-S015-SS Oil Gains Co. Camden PISCO S015-SS      S      0
#:> 4 PSCO-S012-LS Oil Gains Co. <NA> PISCO S012-LS      S      0
#:> 5 PSCO-M001-TS Oil Gains Co. Aida PISCO M001-TS    <NA>      0
#:> 6 PSCO-M0026-TS Oil Gains Co. Ibironke PISCO M026-TS    <NA>      0
#:> # ... with 44 more variables: WellType <dbl>, AL_Method <dbl>,
#:> # Completion <dbl>, SandControl <dbl>, WT_COUNT <dbl>, PVT_GOR <dbl>,
#:> # PVT_API <dbl>, PVT_SG_gas <dbl>, PVT_WaterSalinity <dbl>,
#:> # PVT_H2S <dbl>, PVT_CO2 <dbl>, PVT_PB_CORR <dbl>, PVT_VISC_CORR <dbl>,
#:> # PVT_BPTEMP <dbl>, PVT_BPPRES <dbl>, VLP_CORR <dbl>, IPR_CORR <dbl>,
#:> # IPR_RESPRES <dbl>, IPR_RESTEMP <dbl>, IPR_TOTGOR <dbl>, IPR_WC <dbl>,
#:> # IPR_VOGELRT <dbl>, IPR_VOGELPRES <dbl>, IPR_PI <dbl>, GEO_THMD <chr>,
#:> # GEO_THTEMP <chr>, GL_method <dbl>, GL_ArrayMandrels <chr>,
#:> # GL_Vdepth <dbl>, GL_GSG <dbl>, GL_CO2 <dbl>, WT_DATE <chr>,
#:> # WT_THT <chr>, WT_LIQRT <chr>, WT_WC <chr>, WT_THP <chr>, WT_GOR <chr>,
#:> # WT_GLIR <chr>, WT_DEPTH <chr>, WT_Enable <chr>, WT_GDEPTH <chr>,
#:> # WT_GPRES <chr>, WT_RESPRES <chr>, ProsperFilename <chr>
```

Now it looks much better.

dimensions of the data frame: dim

Let's use some R functions to find out more about our data.

```
# get the dimensions of the table.
dim(myXl)
#:> [1] 100 51
```

Our table has 100 rows and 51 columns.

names of the columns: names

These are the names of the variables or columns:

```
names(myXl)
#:> [1] "Wellname"      "Company"      "Analyst"
#:> [4] "Field"         "Location"     "Platform"
#:> [7] "Fluid"         "WellType"     "AL_Method"
#:> [10] "Completion"    "SandControl"  "WT_COUNT"
#:> [13] "PVT_GOR"       "PVT_API"      "PVT_SG_gas"
#:> [16] "PVT_WaterSalinity" "PVT_H2S"      "PVT_CO2"
#:> [19] "PVT_PB_CORR"   "PVT_VISC_CORR" "PVT_BPTEMP"
#:> [22] "PVT_BPPRES"    "VLP_CORR"     "IPR_CORR"
#:> [25] "IPR_RESPRES"   "IPR_RESTEMP"  "IPR_TOTGOR"
#:> [28] "IPR_WC"        "IPR_VOGELRT"  "IPR_VOGELPRES"
```

```

#> [31] "IPR_PI"          "GEO_THMD"          "GEO_THTEMP"
#> [34] "GL_method"       "GL_ArrayMandrels"  "GL_Vdepth"
#> [37] "GL_GSG"          "GL_CO2"            "WT_DATE"
#> [40] "WT_THT"          "WT_LIQRT"          "WT_WC"
#> [43] "WT_THP"          "WT_GOR"            "WT_GLIR"
#> [46] "WT_DEPTH"        "WT_Enable"         "WT_GDEPTH"
#> [49] "WT_GPRES"        "WT_RESPRES"        "ProsperFilename"

```

summary of the data: summary

```

# A summary of all the variables.
# Notice the difference between numerical and non-numerical variables
summary(myXl)
#>      Wellname      Company      Analyst
#> Length:100      Length:100      Length:100
#> Class :character Class :character Class :character
#> Mode  :character Mode  :character Mode  :character
#>
#>
#>
#>      Field      Location      Platform      Fluid
#> Length:100      Length:100      Length:100      Min.   :0
#> Class :character Class :character Class :character      1st Qu.:0
#> Mode  :character Mode  :character Mode  :character      Median :0
#>                                          Mean    :0
#>                                          3rd Qu.:0
#>                                          Max.    :0
#>
#>      WellType  AL_Method  Completion  SandControl  WT_COUNT
#> Min.   :0      Min.   :0.00  Min.   :0.00  Min.   :0.00  Min.   : 1.00
#> 1st Qu.:0      1st Qu.:1.00  1st Qu.:0.00  1st Qu.:0.00  1st Qu.: 1.00
#> Median :0      Median :1.00  Median :0.00  Median :0.00  Median : 3.00
#> Mean   :0      Mean   :0.98  Mean   :0.07  Mean   :0.24  Mean   : 4.82
#> 3rd Qu.:0      3rd Qu.:1.00  3rd Qu.:0.00  3rd Qu.:0.00  3rd Qu.: 7.00
#> Max.   :0      Max.   :1.00  Max.   :1.00  Max.   :3.00  Max.   :27.00
#>
#>      PVT_GOR      PVT_API      PVT_SG_gas  PVT_WaterSalinity
#> Min.   :280.0      Min.   :35.00  Min.   :0.800  Min.   : 1000
#> 1st Qu.:416.0      1st Qu.:36.00  1st Qu.:1.200  1st Qu.:15000
#> Median :423.0      Median :36.00  Median :1.200  Median :15000
#> Mean   :431.2      Mean   :36.15  Mean   :1.221  Mean   :15247
#> 3rd Qu.:455.2      3rd Qu.:36.00  3rd Qu.:1.237  3rd Qu.:15125
#> Max.   :473.0      Max.   :46.15  Max.   :1.300  Max.   :30000
#>
#>      PVT_H2S      PVT_CO2      PVT_PB_CORR  PVT_VISC_CORR
#> Min.   :0      Min.   :29.00  Min.   :0.00  Min.   :0.00
#> 1st Qu.:0      1st Qu.:65.00  1st Qu.:3.00  1st Qu.:1.00
#> Median :0      Median :65.00  Median :3.00  Median :2.00
#> Mean   :0      Mean   :66.58  Mean   :2.78  Mean   :1.77
#> 3rd Qu.:0      3rd Qu.:69.25  3rd Qu.:3.00  3rd Qu.:2.00
#> Max.   :0      Max.   :74.28  Max.   :3.00  Max.   :4.00

```

```

#:>
#:>      PVT_BPTEMP      PVT_BPPRES      VLP_CORR      IPR_CORR
#:> Min.      : 97.78   Min.      :1683   Min.      : 0.00   Min.      :0.00
#:> 1st Qu.:208.00   1st Qu.:1722   1st Qu.:10.00   1st Qu.:0.00
#:> Median :209.00   Median :1753   Median :10.00   Median :1.00
#:> Mean    :215.58   Mean    :1779   Mean    :10.07   Mean    :0.53
#:> 3rd Qu.:209.00   3rd Qu.:1836   3rd Qu.:10.00   3rd Qu.:1.00
#:> Max.     :408.20   Max.     :1936   Max.     :18.00   Max.     :3.00
#:> NA's     :1       NA's     :1
#:>      IPR_RESPRES      IPR_RESTEMP      IPR_TOTGOR      IPR_WC
#:> Min.      : 658   Min.      :107.6   Min.      : 404   Min.      : 0.00
#:> 1st Qu.:1246   1st Qu.:206.0   1st Qu.: 595   1st Qu.:51.75
#:> Median :1404   Median :209.0   Median : 1247   Median :70.00
#:> Mean    :1386   Mean    :207.3   Mean    : 2028   Mean    :64.73
#:> 3rd Qu.:1565   3rd Qu.:211.0   3rd Qu.: 2348   3rd Qu.:87.53
#:> Max.     :2727   Max.     :226.0   Max.     :11229   Max.     :96.00
#:>
#:>      IPR_VOGELRT      IPR_VOGELPRES      IPR_PI      GEO_THMD
#:> Min.      : 0.0   Min.      : 0.0   Min.      : 0.0000   Length:100
#:> 1st Qu.: 0.0   1st Qu.: 0.0   1st Qu.: 0.8261   Class :character
#:> Median : 559.4   Median : 782.2   Median : 1.7362   Mode  :character
#:> Mean    : 670.3   Mean    : 659.0   Mean    : 2.6829
#:> 3rd Qu.:1145.0   3rd Qu.: 982.0   3rd Qu.: 3.4625
#:> Max.     :2420.8   Max.     :1381.1   Max.     :12.0000
#:>
#:>      GEO_THTEMP      GL_method      GL_ArrayMandrels      GL_Vdepth
#:> Length:100      Min.      :0.00   Length:100      Min.      : 0
#:> Class :character   1st Qu.:0.00   Class :character   1st Qu.:1220
#:> Mode  :character   Median :2.00   Mode  :character   Median :1601
#:>                      Mean    :1.06                      Mean    :2143
#:>                      3rd Qu.:2.00                      3rd Qu.:2304
#:>                      Max.     :2.00                      Max.     :8852
#:>
#:>      GL_GSG      GL_CO2      WT_DATE      WT_THT
#:> Min.      :0.800   Min.      :65.0   Length:100      Length:100
#:> 1st Qu.:1.200   1st Qu.:65.0   Class :character   Class :character
#:> Median :1.200   Median :65.0   Mode  :character   Mode  :character
#:> Mean    :1.196   Mean    :65.1
#:> 3rd Qu.:1.200   3rd Qu.:65.0
#:> Max.     :1.200   Max.     :70.0
#:>
#:>      WT_LIQRT      WT_WC      WT_THP
#:> Length:100      Length:100      Length:100
#:> Class :character   Class :character   Class :character
#:> Mode  :character   Mode  :character   Mode  :character
#:>
#:>
#:>
#:>
#:>      WT_GOR      WT_GLIR      WT_DEPTH
#:> Length:100      Length:100      Length:100
#:> Class :character   Class :character   Class :character
#:> Mode  :character   Mode  :character   Mode  :character

```

```

#:>
#:>
#:>
#:>
#:>   WT_Enable      WT_GDEPTH      WT_GPRES
#:> Length:100      Length:100      Length:100
#:> Class :character Class :character Class :character
#:> Mode  :character Mode  :character Mode  :character
#:>
#:>
#:>
#:>   WT_RESPRES      ProsperFilename
#:> Length:100      Length:100
#:> Class :character Class :character
#:> Mode  :character Mode  :character
#:>
#:>
#:>
#:>

```

structure of the data: **str**

```

# show the data type structure of the table
str(myX1)
#:> Classes 'tbl_df', 'tbl' and 'data.frame': 100 obs. of 51 variables:
#:> $ Wellname      : chr "PSCO-M005-TS" "PSCO-M0007-TS" "PSCO-M004-LS" "PSCO-M008-TS" ...
#:> $ Company       : chr "Oil Gains Co." "Oil Gains Co." "Oil Gains Co." "Oil Gains Co." ...
#:> $ Analyst       : chr "Aida" "Aida" "Aida" "Aida" ...
#:> $ Field        : chr "PISCO" "PISCO" "PISCO" "PISCO" ...
#:> $ Location      : chr "M005-TS" "M007-TS" "M004-LS" "M008-TS" ...
#:> $ Platform      : chr "M" "M" "M" "M" ...
#:> $ Fluid         : num 0 0 0 0 0 0 0 0 0 0 ...
#:> $ WellType      : num 0 0 0 0 0 0 0 0 0 0 ...
#:> $ AL_Method     : num 1 1 1 1 1 0 1 1 1 1 ...
#:> $ Completion    : num 0 0 0 0 0 0 0 0 0 0 ...
#:> $ SandControl   : num 0 0 0 0 0 0 0 0 0 0 ...
#:> $ WT_COUNT      : num 27 22 11 14 13 20 3 2 2 2 ...
#:> $ PVT_GOR       : num 446 473 280 414 420 ...
#:> $ PVT_API       : num 36 36 36 36 35.2 ...
#:> $ PVT_SG_gas    : num 1.2 1.2 1.2 1.2 1.2 ...
#:> $ PVT_WaterSalinity: num 25000 25000 25000 25000 25000 25000 15000 15000 15000 15000 ...
#:> $ PVT_H2S       : num 0 0 0 0 0 0 0 0 0 0 ...
#:> $ PVT_CO2       : num 65.5 65 65 65 65 65 65 65 65 65 ...
#:> $ PVT_PB_CORR   : num 3 3 3 3 3 3 3 3 3 3 ...
#:> $ PVT_VISC_CORR : num 2 2 2 2 2 2 4 2 0 2 ...
#:> $ PVT_BPTEMP    : num 209 209 209 209 209 209 209 209 209 209 ...
#:> $ PVT_BPPRES    : num 1821 1921 1753 1698 1722 ...
#:> $ VLP_CORR      : num 10 10 1 10 10 10 10 10 10 10 ...
#:> $ IPR_CORR      : num 0 0 1 1 0 0 0 1 1 0 ...
#:> $ IPR_RESPRES   : num 930 1300 1573 1286 1468 ...
#:> $ IPR_RESTEMP   : num 209 209 209 200 209 209 214 211 202 216 ...

```



```

#:>      Location      Platform      GEO_THMD      GEO_THTEMP
#:>      "character"    "character"    "character"    "character"
#:>  GL_ArrayMandrels    WT_DATE      WT_THT      WT_LIQRT
#:>      "character"    "character"    "character"    "character"
#:>      WT_WC      WT_THP      WT_GOR      WT_GLIR
#:>      "character"    "character"    "character"    "character"
#:>      WT_DEPTH    WT_Enable    WT_GDEPTH    WT_GPRES
#:>      "character"    "character"    "character"    "character"
#:>      WT_RESPRES  ProsperFilename  Fluid      WellType
#:>      "character"    "character"    "double"    "double"
#:>      AL_Method    Completion    SandControl    WT_COUNT
#:>      "double"      "double"      "double"      "double"
#:>      PVT_GOR      PVT_API      PVT_SG_gas  PVT_WaterSalinity
#:>      "double"      "double"      "double"      "double"
#:>      PVT_H2S      PVT_CO2      PVT_PB_CORR  PVT_VISC_CORR
#:>      "double"      "double"      "double"      "double"
#:>      PVT_BPTEMP    PVT_BPPRES    VLP_CORR      IPR_CORR
#:>      "double"      "double"      "double"      "double"
#:>      IPR_RESPRES    IPR_RESTEMP    IPR_TOTGOR      IPR_WC
#:>      "double"      "double"      "double"      "double"
#:>      IPR_VOGELRT    IPR_VOGELPRES    IPR_PI      GL_method
#:>      "double"      "double"      "double"      "double"
#:>      GL_Vdepth      GL_GSG      GL_CO2
#:>      "double"      "double"      "double"

```

An inventory of the kind of data we have: `table`

```

table(dataTypes)
#:> dataTypes
#:> character      double
#:>      22      29

```