# Attributes and Data Types

**Prof Farhi Marir**

Director of Big Data Analytics Research Lab,
College of Technological Innovation,
Zayed University

# Lecture Content

- Attributes & Data Types,
    - NOIR Attributes,
    - Vectors
    - Arrays and Matrices
    - Data Frames
    - List
    - Descriptive Statistics

# Attributes & Data Types

# NOIR Attributes Types

- NOIR stands for Nominal, Ordinal, Integer and Ration
- Below are the chrematistic of each of them

| | Categorical (Qualitative) | | Numeric (Quantitative) | |
|---|---|---|---|---|
| | **Nominal** | **Ordinal** | **Interval** | **Ratio** |
| **Definition** | The values represent labels that distinguish one from another | Attributes that imply sequence | The difference between two values is meaningful, | Both the difference and the ratio of two values is meaningful |
| **Example** | ZIP code, nationality, street names, employee ID numbers, TRUE or FALSE | Quality of Diamond, academic grades, magnitude of earthquakes, etc. | Temperature in Celsius, Fahrenheit, Calendar dates, latitudes | Age, temperature in Kelvin, counts, length, weight |
| **Associated Operations** | =, ≠ | =, ≠ <br> < , ≤, > , ≥ | =, ≠ <br> < , ≤, > , ≥ <br> +, − | =, ≠ <br> < , ≤, > , ≥ <br> +, −, <br> ×, ÷ |

# Numeric, Character, and Logical Data Types

- R supports numeric, character, and Logical (Boolean) data types
- Example of such code

```
66  # Numeric, Character, and Logical Data Types
67  i <- 1                          # create a numeric variable
68  sport <- "football"             # create a character variable
69  flag <- TRUE                    # create a logical variable
70
71  class(i)                        # returns "numeric"
72  typeof(i)                       # returns "double"
73  class(sport)                    # returns "character"
74  typeof(sport)                   # returns "character"
75  class(flag)                     # returns "logical"
76  typeof(flag)                    # returns "logical"
77  |
78  is.integer(i)                   # returns FALSE
79  j <- as.integer(i)              # coerces contents of i into an integer
80  is.integer(j)                   # returns TRUE
81
```

77:1    (Top Level) ÷                                                        R Script ÷

**Console** c:/data/

```
> flag <- TRUE                     # create a logical variable
> flag <- TRUE                     # create a logical variable
> i <- 1                           # create a numeric variable
> sport <- "football"              # create a character variable
> flag <- TRUE                     # create a logical variable
>
> class(i)                         # returns "numeric"
[1] "numeric"
> typeof(i)                        # returns "double"
[1] "double"
> class(sport)                     # returns "character"
[1] "character"
> typeof(sport)                    # returns "character"
[1] "character"
> class(flag)                      # returns "logical"
[1] "logical"
> typeof(flag)                     # returns "logical"
[1] "logical"
> |
```

# Coerce a variable into a specific type

*is.integer()* *function* can check if a variable is integer or not?

as.integer() can coerce content of a variable into integer (see R code highlights below)

```
77
78    is.integer(i)              # returns FALSE
79    j <- as.integer(i)         # coerces contents of i into an integer
80    is.integer(j)              # returns TRUE
81
80:43   (Top Level)                                                    R Script

Console c:/data/
>
> class(i)                  # returns "numeric"
[1] "numeric"
> typeof(i)                 # returns "double"
[1] "double"
> class(sport)             # returns "character"
[1] "character"
> typeof(sport)            # returns "character"
[1] "character"
> class(flag)              # returns "logical"
[1] "logical"
> typeof(flag)             # returns "logical"
[1] "logical"
> is.integer(i)            # returns FALSE
[1] FALSE
> j <- as.integer(i)       # coerces contents of i into an integer
> is.integer(j)            # returns TRUE
[1] TRUE
>
```

# length() function

length() function reveals the length of a variable. In the example below even Sport ("Football") variable has a length of 1 because it is an element of a vector

# Attributes & Data Types

VECTORS

# Vectors

- Vectors are basic building block for data in R
- A vector can only consist of values in the same class e.g. Vector days (Mon, Tue, Wed, Thu, Fri, Sat, Sun)

```
87  # Vectors
88  is.vector(i)                  # returns TRUE
89  is.vector(flag)               # returns TRUE
90  is.vector(sport)              # returns TRUE
91
92  u <- c("red", "yellow", "blue") # create a vector "red" "yellow" "blue"
93  u                             # returns "red" "yellow" "blue"
94  u[1]                          # returns "red" (1st element in u)
95  v <- 1:5                      # create a vector 1 2 3 4 5
96  v                             # returns 1 2 3 4 5
97  sum(v)                        # returns 15
90:43    (Top Level)                                                     R Script
```

```
Console c:/data/
> is.integer(i)              # returns FALSE
[1] FALSE
> j <- as.integer(i)         # coerces contents of i into an integer
> is.integer(j)              # returns TRUE
[1] TRUE
> length(i)                  # returns 1
[1] 1
> length(flag)               # returns 1
[1] 1
> length(sport)              # returns 1 (not 8 for "football")
[1] 1
> # Vectors
> is.vector(i)               # returns TRUE
[1] TRUE
> is.vector(flag)            # returns TRUE
[1] TRUE
> is.vector(sport)           # returns TRUE
[1] TRUE
>
```

9

# Create a vector and operate on it

Example below shows creation and operations on vectors data type

```
 92   u <- c("red", "yellow", "blue")   # create a vector "red" "yellow" "blue"
 93   u                                  # returns "red" "yellow" "blue"
 94   u[1]                               # returns "red" (1st element in u)
 95   v <- 1:5                           # create a vector 1 2 3 4 5
 96   v                                  # returns 1 2 3 4 5
 97   sum(v)                             # returns 15
 98   w <- v * 2                         # create a vector 2 4 6 8 10
 99   w                                  # returns 2 4 6 8 10
100   w[3]                               # returns 6 (the 3rd element of w)
101   z <- v + w                         # sums two vectors element by element
102   z                                  # returns 3 6 9 12 15
103   z > 8                              # returns FALSE FALSE TRUE TRUE TRUE
104   z[z > 8]                           # returns 9 12 15
105   z[z > 8 | z < 5]                   # returns 3 9 12 15 ("|" denotes "or")
106:1   (Top Level) ⇕                                                    R Script ⇕
```

```
Console c:/data/
[1] 1
> # Vectors
> is.vector(i)            # returns TRUE
[1] TRUE
> is.vector(flag)        # returns TRUE
[1] TRUE
> is.vector(sport)       # returns TRUE
[1] TRUE
> u <- c("red", "yellow", "blue")  # create a vector "red" "yellow" "blue"
> u                       # returns "red" "yellow" "blue"
[1] "red"     "yellow" "blue"
> u[1]                    # returns "red" (1st element in u)
[1] "red"
> v <- 1:5                # create a vector 1 2 3 4 5
> v                       # returns 1 2 3 4 5
[1] 1 2 3 4 5
> sum(v)                  # returns 15
[1] 15
> w <- v * 2              # create a vector 2 4 6 8 10
> w                       # returns 2 4 6 8 10
[1]   2   4   6   8 10
```

```
> w[3]                    # returns 6 (the 3rd element of w)
[1] 6
> z <- v + w              # sums two vectors element by element
> z                       # returns 3 6 9 12 15
[1]   3   6   9 12 15
> z > 8                   # returns FALSE FALSE TRUE TRUE TRUE
[1] FALSE FALSE  TRUE  TRUE  TRUE
> z[z > 8]                # returns 9 12 15
[1]  9 12 15
> z[z > 8 | z < 5]        # returns 3 9 12 15 ("|" denotes "or")
[1]   3   9 12 15
>
```

# Create a Vector with fixed length

- *Vector(length="value")* function create a logical vector in which length could be fixed like in the example below
- You can add new values in the vector
- A vector can be of a different type by using *mode* parameter as shown below

```
107  a <- vector(length=3)              # create a logical vector of length 3
108  a                                  # returns FALSE FALSE FALSE
109  b <- vector(mode="numeric", 3)     # create a numeric vector of length 3
110  typeof(b)                          # returns "double"
111  b[2] <- 3.1                        # assign 3.1 to the 2nd element
112  b                                  # returns 0.0 3.1 0.0
113  c <- vector(mode="integer", 0)     # create an integer vector of length 0
114  c                                  # returns integer(0)
115  length(c)                          # returns 0
116  length(b)                          # returns 3
```
115:44   (Top Level)                                                          R Script

```
Console c:/data/
> z[z > 8 | z < 5]                      # returns 3 9 12 15 ( | denotes  or )
[1]  3  9 12 15
> a <- vector(length=3)                 # create a logical vector of length 3
> a                                     # returns FALSE FALSE FALSE
[1] FALSE FALSE FALSE
> b <- vector(mode="numeric", 3)        # create a numeric vector of length 3
> typeof(b)                             # returns "double"
[1] "double"
> b[2] <- 3.1                           # assign 3.1 to the 2nd element
> b                                     # returns 0.0 3.1 0.0
[1] 0.0 3.1 0.0
> c <- vector(mode="integer", 0)        # create an integer vector of length 0
> c                                     # returns integer(0)
integer(0)
> length(c)                             # returns 0
```

# Attributes & Data Types

ARRAYS and MATRICES

# Arrays

arrays() function can be used to restructure a vector as an array

Highlighted R code below builds a three dimensional array to hold the quarterly sales for three region over two-year period

```
119  # Arrays and Matrices
120
121  # the dimensions are 3 regions, 4 quarters, and 2 years
122  quarterly_sales <- array(0, dim=c(3,4,2))
123  quarterly_sales[2,1,1] <- 158000
124  quarterly_sales
125
126  sales_matrix <- matrix(0, nrow = 3, ncol = 4)
```
124:16    (Top Level) ⇕                                                    R Script ⇕

Console c:/data/  ⇗
> dim(b)                        # returns NULL (an undefined value)
NULL
> # the dimensions are 3 regions, 4 quarters, and 2 years
> quarterly_sales <- array(0, dim=c(3,4,2))
> quarterly_sales[2,1,1] <- 158000
> quarterly_sales
, , 1

      [,1] [,2] [,3] [,4]
[1,]     0    0    0    0
[2,] 158000    0    0    0
[3,]     0    0    0    0

, , 2

      [,1] [,2] [,3] [,4]
[1,]     0    0    0    0
[2,]     0    0    0    0
[3,]     0    0    0    0

# Matrix

- A two dimension array is known as a **matrix**
- The following R code initializes a matrix to hold the quarterly sale for three regions. It uses *matrix()* function with *nrow* & *ncol* as parameters to define number of rows and columns respectively for *Sales_matrix*

```
126  sales_matrix <- matrix(0, nrow = 3, ncol = 4)
127  sales_matrix
128
129  install.packages("matrixcalc")          # install, if necessary
127:13   (Top Level)                                                        R Script
```

```
Console c:/data/
, , 1

        [,1] [,2] [,3] [,4]
[1,]       0    0    0    0
[2,] 158000    0    0    0
[3,]       0    0    0    0

, , 2

        [,1] [,2] [,3] [,4]
[1,]     0    0    0    0
[2,]     0    0    0    0
[3,]     0    0    0    0

> sales_matrix <- matrix(0, nrow = 3, ncol = 4)
> sales_matrix
        [,1] [,2] [,3] [,4]
[1,]     0    0    0    0
[2,]     0    0    0    0
[3,]     0    0    0    0
```

# Operations on Matrix

- R provides the standard matrix operations such as addition, subtraction, and multiplication as well the transpose function t() and the inverse matrix function *matrix.inverse()*

- For this to run you need to include *the matrixcalc* package
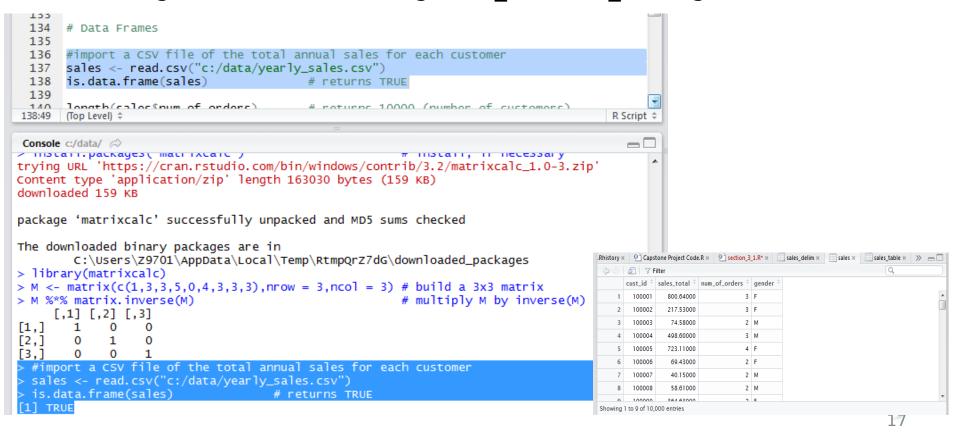
```
129  install.packages("matrixcalc")                    # install, if necessary
130  library(matrixcalc)
131  M <- matrix(c(1,3,3,5,0,4,3,3,3),nrow = 3,ncol = 3) # build a 3x3 matrix
132  M %*% matrix.inverse(M)                            # multiply M by inverse(M)
133
134  # Data Frames
```

132:79   (Top Level) ÷                                                      R Script ÷

Console c:/data/

```
     [,1] [,2] [,3] [,4]
[1,]    0    0    0    0
[2,]    0    0    0    0
[3,]    0    0    0    0
> install.packages("matrixcalc")                         # install, if necessary
trying URL 'https://cran.rstudio.com/bin/windows/contrib/3.2/matrixcalc_1.0-3.zip'
Content type 'application/zip' length 163030 bytes (159 KB)
downloaded 159 KB

package 'matrixcalc' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
        C:\Users\Z9701\AppData\Local\Temp\RtmpQrZ7dG\downloaded_packages
> library(matrixcalc)
> M <- matrix(c(1,3,3,5,0,4,3,3,3),nrow = 3,ncol = 3) # build a 3x3 matrix
> M %*% matrix.inverse(M)                            # multiply M by inverse(M)
     [,1] [,2] [,3]
[1,]    1    0    0
[2,]    0    1    0
[3,]    0    0    1
```

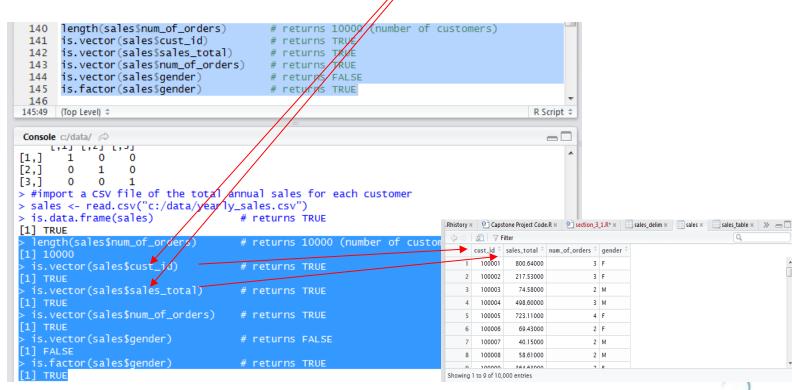# Attributes and Data Types

Data Frames

# Data Frames

- Similar to the concepts, data frames provide a structure for storing and accessing several variables of different data types,

- In fact *read.csv()* in the accompanying R example created data frames containing different variable e.g. cust_ID, sales_total, gender, etc..
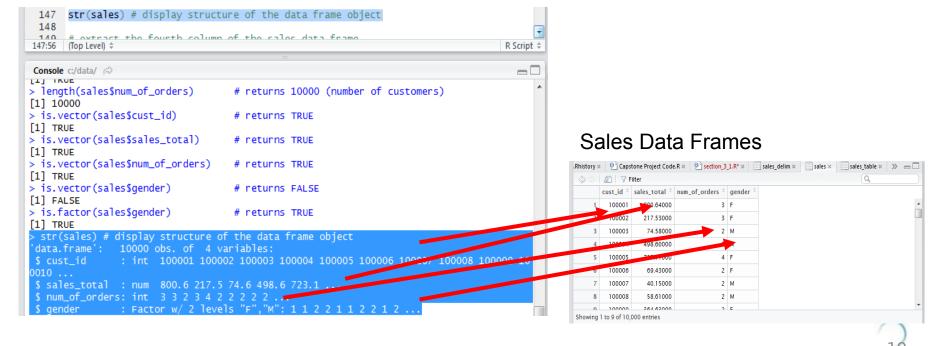
# Access through $ notation

- The data stored in data frame can be accessed using $ notation,

# str() function

- Because of their flexibility to handle many data types, are the preferred input format for many of the modelling function,

- *str()* function display the structure & values of the data frame e.g. *Sales below*,

- This function identifies integer and numeric (double) data types, the factor variable and levels, as well as first value of each variable



Sales Data Frames

# Retrieve values/data from data frames

- Data frames are lists of variables of the same length,
- Subset of the data frame can be retrieved from **sub-setting operators**
- R's sub-setting operators are powerful as they allow complex operations to retrieve subset of the data set

sales[,4]: Extract the fourth Column of Sales

**sales$gender** retrieves gender column as shown below

# Retrieve rows

*sales[1:2,]* retrieves the first two rows

# Retrieve selected Columns

- *sales [,c(1,2,3)]* retrieves 1, 2,  and 3<sup>rd</sup> columns

# Retrieve named Columns

*sales[,c("cust_id", "sales_total")]* retrieves cust_id and sales_total columns

```
157   # retrieve both the cust_id and the sales_total columns
158   sales[,c("cust_id", "sales_total")]
159   # retrieve all the records whose gender is female
```
158:36   (Top Level) ÷                                                    R Script ÷

**Console** c:/data/

```
4002    104002      46.36000
4003    104003     147.94000
4004    104004      41.44000
4005    104005      66.76000
4006    104006      63.58000
4007    104007      80.05000
4008    104008     128.03000
4009    104009     198.38000
4010    104010     100.96000
4011    104011      40.82000
4012    104012     693.41000
4013    104013     230.96000
4014    104014      69.66000
4015    104015    1558.03000
4016    104016      94.47000
4017    104017     136.53000
4018    104018      33.23000
4019    104019     208.12000
4020    104020     236.27000
4021    104021    1086.78000
4022    104022     557.65000
```

# Retrieve Columns based on a predicate

*sales[sales$gender=="F",]* retrieves records whose gender is female i.e. "F"

```
158   sales[,c("cust_id", "sales_total")]
159   # retrieve all the records whose gender is female
160   sales[sales$gender=="F",]
161
162   class(sales)
```

160:26   (Top Level) ⇕                                                              R Script ⇕

Console c:/data/ ⇗

```
3071   103071      82.41000              1        F
3073   103073     132.36000              2        F
3075   103075   6428.06000             20        F
3076   103076     481.97000              3        F
3079   103079     225.31000              2        F
3080   103080     353.02000              3        F
3081   103081     262.37000              4        F
3082   103082     113.28000              1        F
3085   103085      81.77000              1        F
3091   103091     133.23000              3        F
3092   103092     212.32000              2        F
3093   103093     232.55000              3        F
3095   103095      49.46600              2        F
3097   103097     149.12000              1        F
3099   103099     156.95000              2        F
3101   103101      56.80000              2        F
3103   103103     256.48000              1        F
3104   103104      67.73000              3        F
3106   103106     531.52000              5        F
3107   103107      49.31000              1        F
3108   103108     117.05000              1        F
```

Attributes & Data Type here

**LISTS**

# Lists

- A list is a collection of objects that can be of various type, including other lists
- *list()* function is used to create lists
- Using the Vector **v** and the Matrix *M* created in earlier example, the following R code creates assortments which is a list of different object types

# Displaying Content of a List

- Single bracket [ ] in the list() function only accesses an item not its content,
- Double bracket [[ ]] instead is used to access the content
- See highlighted R code and results below where [ ] and [[ ]] are used giving items and content respectively

```
175  # examine the fifth object, M, in the list
176  class(assortment[5])           # returns "list"
177  length(assortment[5])          # returns 1
178  class(assortment[[5]])         # returns "matrix"
179  length(assortment[[5]])        # returns 9 (for the 3x3 matrix)
180
180:1   (Top Level)                                              R Script
```

Console C:/Users/Z9701/Desktop/aaTeaching/CTI Courses/Fall 2016/CTI 466 Data Analytics/CIT466 R Exercises/

```
[[4]]
[1] 1 2 3 4 5

[[5]]
     [,1] [,2] [,3]
[1,]   1    5    3
[2,]   3    0    3
[3,]   3    4    3

> # examine the fifth object, M, in the list
> class(assortment[5])          # returns "list"
[1] "list"
> length(assortment[5])         # returns 1
[1] 1
> class(assortment[[5]])        # returns "matrix"
[1] "matrix"
> length(assortment[[5]])       # returns 9 (for the 3x3 matrix)
[1] 9
```

# Use of str() function in Lists

- *str()* function display details of the structure of the list

# Descriptive Statistics

- It has already been shown that the summary() function provides several descriptive statistics, such as the mean and median, about a variable such as the *sales* data frame.

- The following code provides some common R functions that include descriptive statistics. In parentheses, the comments describe the functions.

# Descriptive Statistics

```
# to simplify the function calls, assign
x <- sales$sales_total
y <- sales$num_of_orders
cor(x,y) # returns 0.7508015 (correlation)
cov(x,y) # returns 345.2111 (covariance)
IQR(x) # returns 215.21 (interquartile range)
mean(x) # returns 249.4557 (mean)
median(x) # returns 151.65 (median)
range(x) # returns 30.02 7606.09 (min max)
sd(x) # returns 319.0508 (std. dev.)
var(x) # returns 101793.4 (variance)
```

# **Descriptive Statistics**

- The function apply() is useful when the same function is to be applied to several variables in a data frame. For example, the following R code calculates the standard deviation for the first three variables in *sales.*

- apply(sales[,c(1:3)], MARGIN=2, FUN=sd)

```
        cust_id    sales_total  num_of_orders
   2886.895680    319.050782        1.441119
```