DATA SCIENCE WITH R



Class 4 – Data Manipulation in R

Topic 1



Manipulating Data Using Base R



INDEX



Manipulating data using base R

Using dplyr to manipulate data

Working with date objects

Merging tables

Missing value treatment

Using reshape2() to transpose data

Manipulating Character Strings

Using sqldf

Data Manipulation: Base R

Data Manipulation: Base R

- Sub-setting data
- Selecting specified columns
- Adding new columns
- Reordering data (Ascending/Descending order)
- Group wise operations
- Producing contingency tables

Sub-setting data

Sub setting: Selecting a sub set of rows across all columns

```
> head(oj[oj$brand=='tropicana',])
            brand week logmove feat price
                                               AGE 60
  store
      2 tropicana
                    40 9.018695
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
     2 tropicana
                   46 8.723231
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
     2 tropicana
                    47 8.253228
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
     2 tropicana
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
                    48 8.987197
     2 tropicana
                    50 9.093357
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
     2 tropicana
                    51 8.877382
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
    HVAL150 SSTRDIST SSTRVOL CPDIST5
1 0.4638871 2.110122 1.142857 1.92728 0.3769266
2 0.4638871 2.110122 1.142857 1.92728 0.3769266
3 0.4638871 2.110122 1.142857 1.92728 0.3769266
4 0.4638871 2.110122 1.142857 1.92728 0.3769266
5 0.4638871 2.110122 1.142857 1.92728 0.3769266
6 0.4638871 2.110122 1.142857 1.92728 0.3769266
```

Can use multiple conditions, | (or), & (and) operator

```
> head(oj[oj$brand=='tropicana'|oj$brand=='dominicks',])
            brand week logmove feat price
  store
                                                AGE 60
                                                           EDUC
                                                                            INCOME
      2 tropicana
                    40 9.018695
                                       3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
      2 tropicana
                    46 8.723231
                                       3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
      2 tropicana
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
                    47 8.253228
      2 tropicana
                    48 8.987197
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
      2 tropicana
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
                    50 9.093357
      2 tropicana
                    51 8.877382
                                       3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
                               SSTRVOL CPDIST5
    WORKWOM
              HVAL150 SSTRDIST
                                                   CPWVOL 5
1 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266
2 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266
3 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266
4 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266
5 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266
6 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266
> dim(oj[oj$brand=='tropicana'|oj$brand=='dominicks',])
[1] 19298
```

```
> dim(oj[oj$brand=='tropicana' & oj$feat==0,])
[1] 8045
> head(oj[oj$brand=='tropicana' & oj$feat==0,])
            brand week logmove feat price
                                                AGE 60
                                                           EDUC
                                                                   ETHNIC
                                                                             INCOME
                                                                                      HHLARGE
      2 tropicana
                    40 9.018695
                                       3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
      2 tropicana
                    46 8.723231
                                       3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
      2 tropicana
                    47 8.253228
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
      2 tropicana
                    48 8.987197
                                       3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
                    50 9.093357
      2 tropicana
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
      2 tropicana
                    51 8.877382
                                       3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
              HVAL150 SSTRDIST
                                SSTRVOL CPDIST5
    WORKWOM
                                                   CPWVOL5
1 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266
2 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266
3 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266
 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266
5 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266
 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266
```

- So, far logical sub-setting is discussed.
- Use which() operator to get the index for specific rows

```
> index<-which(oj$brand=="dominicks")
> head(index)
[1] 221 222 223 224 225 226
> head(oi[index,])
                           logmove feat price
    store
              brand week
        2 dominicks
                                         1.59 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853 0.4638871
221
                          9.264829
222
        2 dominicks
                      46 8.987197
                                         2.69 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853 0.4638871
223
        2 dominicks
                          8.831712
                                         2.09 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853 0.4638871
224
        2 dominicks
                      48 7.965546
                                         2.09 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853 0.4638871
225
        2 dominicks
                      50 7.377759
                                         2.09 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853 0.4638871
                                         1.89 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853 0.4638871
        2 dominicks
    SSTRDIST SSTRVOL CPDIST5
221 2.110122 1.142857 1.92728 0.3769266
222 2.110122 1.142857 1.92728 0.3769266
223 2.110122 1.142857 1.92728 0.3769266
224 2.110122 1.142857 1.92728 0.3769266
225 2.110122 1.142857 1.92728 0.3769266
226 2.110122 1.142857 1.92728 0.3769266
```

Logical vectors Vs. which

- which() removes NA values in the logical vector
- It only returns the indices where the logical vector is TRUE

```
> #Consider vector sales with missing values
> sales<-c(100,200,NA,300,400,NA,500,600,700,NA,1000,1500,NA,NA)
> #subset data using logical operator
> sales[sales>600]
[1] NA NA 700 NA 1000 1500 NA NA
> #subset data using which
> sales[which(sales>600)]
[1] 700 1000 1500
```

Selecting Columns

Manipulating data: Base R (Selecting)

Selecting a specified set of columns

Selecting + Sub-setting

```
> head(oj[oj$brand=='tropicana' & oj$feat==0,c("week","store")])
  week store
1    40    2
2    46    2
3    47    2
4    48    2
5    50    2
6    51    2
> dim(oj[oj$brand=='tropicana' & oj$feat==0,c("week","store")])
[1] 8045    2
```

Adding new columns

Adding new columns

```
> oj$logInc<-log(oj$INCOME)
> head(oj)
            brand week logmove feat price
  store
                                                AGE 60
                                                           EDUC
                                                                   ETHNIC
                                                                            INCOME
                                                                                     HHLARGE
      2 tropicana
                    40 9.018695
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
      2 tropicana
                    46 8.723231
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
                    47 8.253228
      2 tropicana
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
      2 tropicana
                    48 8.987197
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
      2 tropicana
                    50 9.093357
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
                    51 8.877382
                                      3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534
      2 tropicana
              HVAL150 SSTRDIST SSTRVOL CPDIST5
                                                           logInc
    WORKWOM
                                                  CPWVOL5
1 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266 2.35643
2 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266 2.35643
3 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266 2.35643
4 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266 2.35643
5 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266 2.35643
6 0.3035853 0.4638871 2.110122 1.142857 1.92728 0.3769266 2.35643
>
```

Ordering data

Ordering

order() returns the element order that results in a sorted vector

```
> students<-c("John","Tim","Alice","Zeus")
> students
[1] "John" "Tim" "Alice" "Zeus"
> order(students)
[1] 3 1 2 4
> students[order(students)]
[1] "Alice" "John" "Tim" "Zeus"
```

Application: Very useful for sorting dataframes

Ordering data

```
> head(oj[order(oj$week),])
                brand week
                            logmove feat price
    store
                                                    AGE 60
                                                               EDUC
                                                                                 INCOME
                        40 9.018695
1
            tropicana
                                       0 3.87 0.2328647 0.2489349 0.11427995 10.55321 0.1039534
111
        2 minute.maid
                        40 8.407378
                                          3.17 0.2328647 0.2489349 0.11427995 10.55321 0.1039534
221
            dominicks
                        40 9.264829
                                       1 1.59 0.2328647 0.2489349 0.11427995 10.55321 0.1039534
331
            tropicana
                        40 8.680672
                                       0 3.66 0.1173680 0.3212257 0.05387528 10.92237 0.1030916
        5 minute.maid
                        40 8.348538
447
                                          2.99 0.1173680 0.3212257 0.05387528 10.92237 0.1030916
563
            dominicks
                        40 7,491088
                                          1.59 0.1173680 0.3212257 0.05387528 10.92237 0.1030916
      WORKWOM
                HVAL150 SSTRDIST
                                   SSTRVOL
                                           CPDIST5
                                                       CPWVOL 5
    0.3035853 0.4638871 2.110122 1.1428571 1.927280 0.3769266
111 0.3035853 0.4638871 2.110122 1.1428571 1.927280 0.3769266
221 0.3035853 0.4638871 2.110122 1.1428571 1.927280 0.3769266
331 0.4105680 0.5358834 3.801998 0.6818182 1.600573 0.7363068
447 0.4105680 0.5358834 3.801998 0.6818182 1.600573 0.7363068
563 0.4105680 0.5358834 3.801998 0.6818182 1.600573 0.7363068
```

Ordering data

```
> head(oj[order(-oj$week),])
                             logmove feat price
    store
                brand week
                                                     AGE 60
                                                                EDUC
            tropicana
                       160
                            8.669743
                                           2.97 0.2328647 0.2489349 0.11427995 10.55321 0.1039534
110
        2 minute.maid
220
                       160 10.626582
                                           2.19 0.2328647 0.2489349 0.11427995 10.55321 0.1039534
330
            dominicks
                           9.064158
                       160
                                           1.82 0.2328647 0.2489349 0.11427995 10.55321 0.1039534
            tropicana
446
                       160 8.921057
                                           2.78 0.1173680 0.3212257 0.05387528 10.92237 0.1030916
        5 minute.maid
                       160 10.825840
                                           2.19 0.1173680 0.3212257 0.05387528 10.92237 0.1030916
562
678
                                            1.85 0.1173680 0.3212257 0.05387528 10.92237 0.1030916
            dominicks
                       160
                            8.723231
                HVAL150 SSTRDIST
                                            CPDIST5
      WORKWOM
                                                       CPWVOL 5
110 0.3035853 0.4638871 2.110122 1.1428571 1.927280 0.3769266
220 0.3035853 0.4638871 2.110122 1.1428571 1.927280 0.3769266
330 0.3035853 0.4638871 2.110122 1.1428571 1.927280 0.3769266
446 0.4105680 0.5358834 3.801998 0.6818182 1.600573 0.7363068
562 0.4105680 0.5358834 3.801998 0.6818182 1.600573 0.7363068
678 0.4105680 0.5358834 3.801998 0.6818182 1.600573 0.7363068
```

- Subsetting data: Using logical subsets and which() statement
- Selecting columns: Using column names at column index
- Adding new columns: Use of \$ operator
- Re-ordering data: order()
- Group Wise Summaries
- Producing Contingency tables

GroupWise operations

- GroupWise operations
- tapply(), aggregate()
- What is the mean price of each brand of juice across all stores?

- GroupWise operations
- tapply(), aggregate()
- What is the mean income level corresponding to brand of juice across all stores?

Contingency tables

Category wise counts: Contingency tables

Income	Age	Gender	Location
10,000,000	24	М	Arizona
20,000,000	32	F	California
15,000,000	28	М	Arizona
18,000,000	26	F	California

Category wise counts: Contingency tables

Counts	California	Arizona
Male	0	2
Female	2	0

Income	California	Arizona
Male	0	10,000,000+15,000,000
Female	20,000,000+18,000,000	

- Category wise counts: Contingency tables
- table(), xtab()
- Number of people who bought different brands categorized by presence of advertising campaigns
 - > table(oj\$brand,oj\$feat)

```
0 1
dominicks 7169 2480
minute.maid 6865 2784
tropicana 8045 1604
```

- Category wise counts: Contingency tables
- table(), xtab()
- Total income categorized by brand and presence of advertisements



RECAP

- Sub-setting data
- Selecting specified columns
- Adding new columns
- Reordering data (Ascending/Descending order)
- Group wise operations
- Producing contingency tables



Topic 2



Using dplyr to Manipulate Data



INDEX



Manipulating data using base R

Using dplyr to manipulate data

Working with date objects

Merging tables

Missing value treatment

Using reshape2() to transpose data

Manipulating Character Strings

Using sqldf

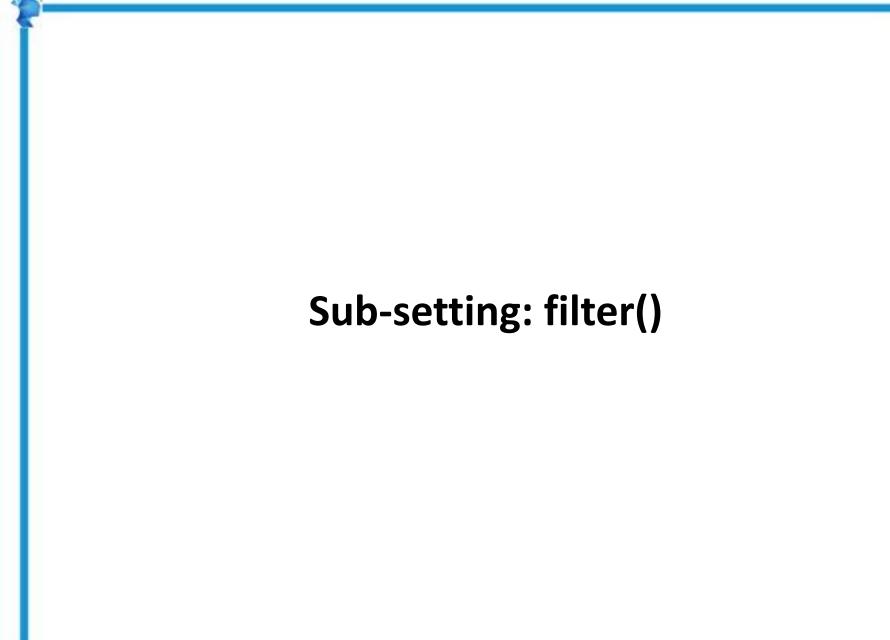


Manipulating data: dplyr

- dplyr: Whats and Whys
- Sub-setting data using filter()
- Selecting columns using select()
- Adding new columns using mutate()
- Ordering data using arrange()
- Summarizing using summarize() and group_by()
- Using functional pipelines to do more than one manipulation task

Manipulating data: dplyr

- Base R: Good for Medium sized data sets, Awkward Syntax
- dplyr: Faster and elegant syntax
- dplyr: Dataframes
- install.packages("dplyr")
- library(dplyr)



- Sub-setting the data using filter(), base R equivalents: logical subsets and which()
- Only that portion of data such that brand bought is "tropicana"

```
> library(dplyr)
> head(filter(oj,brand=="tropicana"))
  store
            brand week logmove feat price
      2 tropicana
                    40 9.018695
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
      2 tropicana
                    46 8.723231
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
      2 tropicana
                    47 8.253228
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
      2 tropicana
                    48 8.987197
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
      2 tropicana
                    50 9.093357
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
      2 tropicana
                    51 8.877382
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
    HVAL150 SSTRDIST SSTRVOL CPDIST5
                                        CPWVOL5
1 0.4638871 2.110122 1.142857 1.92728 0.3769266
2 0.4638871 2.110122 1.142857 1.92728 0.3769266
3 0.4638871 2.110122 1.142857 1.92728 0.3769266
4 0.4638871 2.110122 1.142857 1.92728 0.3769266
5 0.4638871 2.110122 1.142857 1.92728 0.3769266
6 0.4638871 2.110122 1.142857 1.92728 0.3769266
```

- Sub-setting the data using filter(), base R equivalents: logical subsets and which()
- Only that portion of data such that brand bought is "tropicana" or "dominicks"

```
> head(filter(oj,brand=="tropicana"|brand=="dominicks"))
            brand week logmove feat price
  store
                                                           EDUC
                                                                  ETHNIC
      2 tropicana
                    40 9.018695
      2 tropicana
                    46 8.723231
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853 0.4638871 2.110122
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853 0.4638871 2.110122
      2 tropicana
                    47 8.253228
      2 tropicana
                    48 8.987197
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853 0.4638871 2.110122
                    50 9.093357
      2 tropicana
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853 0.4638871 2.110122
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853 0.4638871 2.110122
      2 tropicana
   SSTRVOL CPDIST5
                   CPWVOL 5
   .142857 1.92728 0.3769266
   .142857 1.92728 0.3769266
3 1.142857 1.92728 0.3769266
4 1.142857 1.92728 0.3769266
5 1.142857 1.92728 0.3769266
6 1.142857 1.92728 0.3769266
```





- Selecting columns from data using select(), base R equivalents: index subsets
- Selecting columns brand and income

- Selecting columns from data using select(), base R equivalents: index subsets
- Dropping columns brand and income



- Adding columns to data using mutate(),
- Adding a new column, log(income)

```
> dim(oj)
[1] 28947
> head(mutate(oj,logIncome=log(INCOME)))#Changes not made in oj but its copy
            brand week logmove feat price
  store
      2 tropicana
                    40 9.018695
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
      2 tropicana
                    46 8.723231
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
      2 tropicana
                    47 8.253228
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
      2 tropicana
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
      2 tropicana
                    50 9.093357
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
      2 tropicana
                    51 8.877382
                                   0 3.87 0.2328647 0.2489349 0.1142799 10.55321 0.1039534 0.3035853
    HVAL150 SSTRDIST SSTRVOL CPDIST5
                                        CPWVOL5 logIncome
                                                  2.35643
1 0.4638871 2.110122 1.142857 1.92728 0.3769266
2 0.4638871 2.110122 1.142857 1.92728 0.3769266
                                                  2.35643
3 0.4638871 2.110122 1.142857 1.92728 0.3769266
                                                  2.35643
4 0.4638871 2.110122 1.142857 1.92728 0.3769266
                                                  2.35643
5 0.4638871 2.110122 1.142857 1.92728 0.3769266
                                                  2.35643
6 0.4638871 2.110122 1.142857 1.92728 0.3769266
                                                  2.35643
> dim(oi)
             17
[1] 28947
```

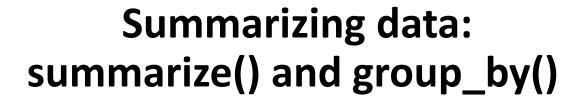


- Ordering data using order_by(),
- Order whole data by income in ascending order

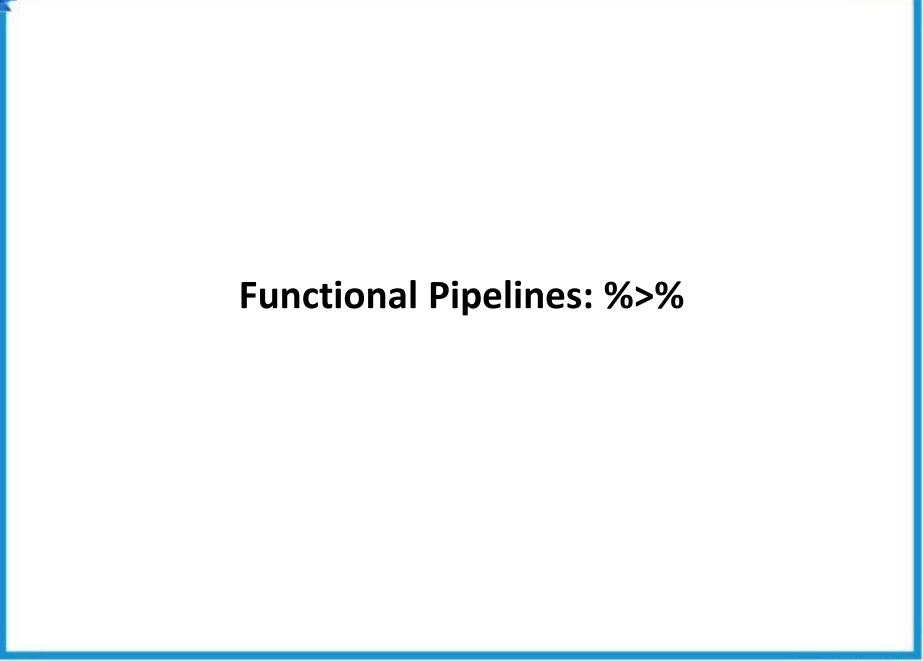
```
> head(arrange(oj,INCOME))
            brand week logmove feat price
  store
                                                AGE 60
                                                           EDUC
                                                                   ETHNIC
                                                                            INCOME
                                                                                      HHLARGE
     75 tropicana
                    40 8.971067
                                   0 3.87 0.2076995 0.2195485 0.4159995 9.867083 0.06396471 0.3155833
    75 tropicana
                    41 8.392990
                                      3.87 0.2076995 0.2195485 0.4159995 9.867083 0.06396471 0.3155833
    75 tropicana
                    42 9.018695
                                      3.87 0.2076995 0.2195485 0.4159995 9.867083 0.06396471 0.3155833
    75 tropicana
                    43 8.624791
                                      3.87 0.2076995 0.2195485 0.4159995 9.867083 0.06396471 0.3155833
    75 tropicana
                    44 8.476371
                                      3.87 0.2076995 0.2195485 0.4159995 9.867083 0.06396471 0.3155833
     75 tropicana
                                      3.87 0.2076995 0.2195485 0.4159995 9.867083 0.06396471 0.3155833
  HVAL150 SSTRDIST
                                       CPWVOL 5
    0.496 7.192667 2.230769 1.375126 0.7031819
    0.496 7.192667 2.230769 1.375126 0.7031819
    0.496 7.192667 2.230769 1.375126 0.7031819
    0.496 7.192667 2.230769 1.375126 0.7031819
    0.496 7.192667 2.230769 1.375126 0.7031819
    0.496 7.192667 2.230769 1.375126 0.7031819
```

- Ordering data using order_by(),
- Order whole data by income in descending order

```
> head(arrange(oj,-INCOME)
            brand week logmove feat price
  store
                                               AGE 60
                                                          EDUC
     62 tropicana
                    40 9.373819
                                   0 3.87 0.2225343 0.5177603 0.0265109 11.2362 0.1039793 0.3227652
     62 tropicana
                    41 9.368369
                                     3.87 0.2225343 0.5177603 0.0265109 11.2362 0.1039793 0.3227652
     62 tropicana
                                     3.87 0.2225343 0.5177603 0.0265109 11.2362 0.1039793 0.3227
                    42 9.570529
    62 tropicana
                                      3.87 0.2225343 0.5177603 0.0265109 11.2362 0.1039793 0.3227652
                    43 9.400630
                                   0 3.87 0.2225343 0.5177603 0.0265109 11.2362 0.1039793 0.3227652
     62 tropicana
                    44 9.329367
     62 tropicana
                                   0 3.87 0.2225343 0.5177603 0.0265109 11.2362 0.1039793 0.3227652
    HVAL150 SSTRDIST
1 0.9166995 5.452685 0.7058824 2.18405 0.2017224
2 0.9166995 5.452685 0.7058824 2.18405 0.2017224
3 0.9166995 5.452685 0.7058824 2.18405 0.2017224
4 0.9166995 5.452685 0.7058824 2.18405 0.2017224
5 0.9166995 5.452685 0.7058824 2.18405 0.2017224
6 0.9166995 5.452685 0.7058824 2.18405 0.2017224
```



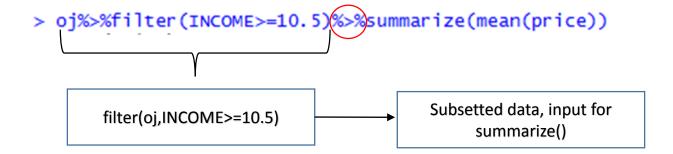
- Summarizing data using summarize() and group_by()
- group_by() makes grouped table, summarize() can take this grouped table and produce summaries for different columns
- Mean level of income and standard deviation of income for each brand of orange juice



- dplyr becomes a powerful tool when combined with %>% (pipe) operator
- Several data manipulation tasks can be accomplished in just one line of code
- Traditionally functional composition is achieved by using nested function calls
- For example, Find the mean price for all people whose income is >=10.5

> oj%>%filter(INCOME>=10.5)%>%summarize(mean(price))

filter(oj,INCOME>=10.5)



Clearly the code looks very messy, using a %>% operator, we can make it more readable

```
> oj%>%filter(INCOME>=10.5)%>%summarize(mean(price))
mean(price)
1 2.270229
```

- This can be easily read as:
- Take data oj, filter it based on income
- Take this filtered data frame and compute the mean of price



 Subset the data based on price>=2.5, create a column logIncome, compute the mean, standard deviation and median of column logIncome

RECAP

- dplyr: better manipulation functionality
- Sub-setting data using filter()
- Selecting columns using select()
- Adding new columns using mutate()
- Ordering data using arrange()
- Summarizing using summarize() and group_by()
- Using functional pipelines to do more than one manipulation task



Topic 3



Working with Date Objects







Using dplyr to manipulate data



Working with date objects

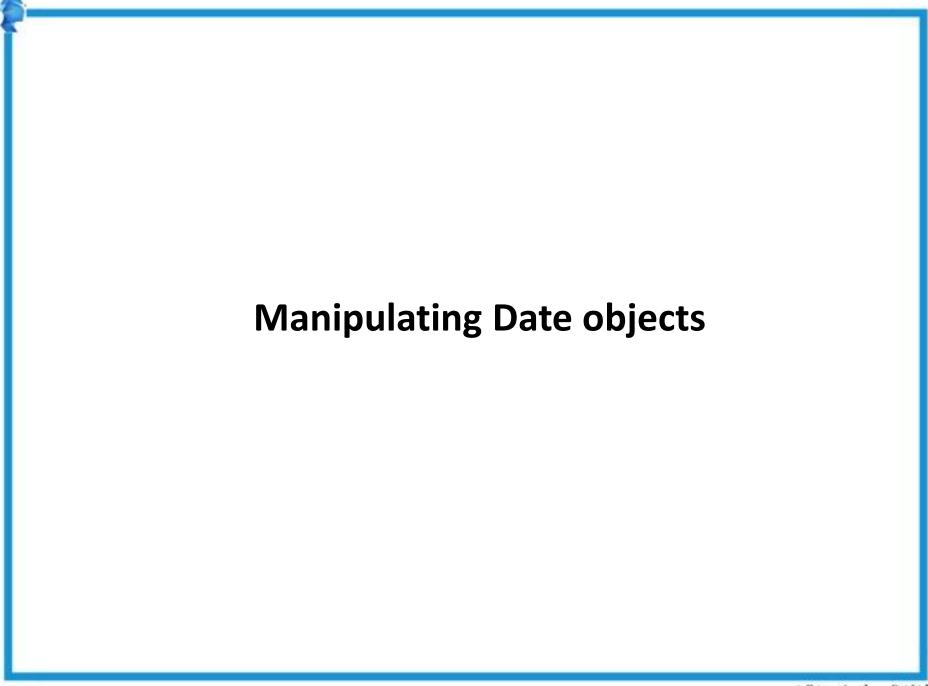
Merging tables

Missing value treatment

Using reshape2() to transpose data

Manipulating Character Strings

Using sqldf



- Dates are treated as a special data type in most programming languages
- R also treats dates as a separate data type
- Doing so, one can easily work with dates
- Usual date operations include:
 - > Finding time interval between two points in data: age
 - > Extracting months and week days

- Character to date conversion-Date Class
- Extracting months and weekdays
- Using difftime()
- Manipulating data involving dates
- POSIXct and POSIXIt Classes
- Working with lubridate()

Using Date class to convert a character into a Date

```
> fd<-read.csv("Fd.csv")
> str(fd)
 'data.frame': 30443 obs. of 25 variables:
   $ FlightDate
                                                       : Factor w/ 74 levels "01-Feb-14", "01-Jan-14",...: 2 2 2 2 2 2 2 2 2 2 ...
   $ UniqueCarrier
                                                       : Factor w/ 13 levels "AA", "AS", "B6", ...: 4 4 4 4 11 11 11 11 11 4 ...
                                                       : int 19790 19790 19790 19790 20355 20355 20355 20355 20355 19790 ...
   $ AirlineID
                                                       : Factor w/ 13 levels "AA", "AS", "B6",...: 4 4 4 4 11 11 11 11 11 14 ...
: Factor w/ 2816 levels "", "D942DN", "NOEGMQ",...: 2641 2512 2490 2581 1657 34 60 17
   $ Carrier
   $ TailNum
98 1443 473 ...
  $ FlightNum
                                                       : int 335 1095 2422 1607 657 894 1843 2041 413 2030 ...
   $ OriginAirportID
                                                       : int 11057 11057 11057 11057 11057 11057 11057 11057 11057 13232 ...
   $ OriginAirportSeqID: int 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 125703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703
```

Using Date class to convert a character into a Date

```
> fd$FlightDate<-as.Date(fd$FlightDate,"%d-%b-%y")</pre>
> str(fd)
'data.frame': 30443 obs. of 25 variables:
 $ FlightDate
                     : Date, format: "2014-01-01" "2014-01-01" "2014-01-01" ...
                     : Factor w/ 13 levels "AA", "AS", "B6", ...: 4 4 4 4 11 11 11 11 11 4 ...
 $ UniqueCarrier
 $ AirlineID
                     : int 19790 19790 19790 19790 20355 20355 20355 20355 20355 19790 ...
                     : Factor w/ 13 levels "AA", "AS", "B6", ...: 4 4 4 4 11 11 11 11 11 4 ...
 $ Carrier
                     : Factor w/ 2816 levels "", "D942DN", "N0EGMQ"...: 2641 2512 2490 2581 1657 34 60 17
 $ TailNum
98 1443 473 ...
                     : int 335 1095 2422 1607 657 894 1843 2041 413 2030 ...
 $ FlightNum
                     : int 11057 11057 11057 11057 11057 11057 11057 11057 11057 13232 ...
 $ OriginAirportID
$ OriginAirportSeqID: int 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 132
```



Using Date class to convert a character into a Date

Code	Value
%d	Day of month (decimal number)
%m	Month (decimal number)
%b	Month (abbreviated)
%B	Month (full name)
%у	Year (2 digits)
%Y	Year (4 digits)

25/Aug/04: "%d/%b/%y

25-August-2004:%d-%B-%Y



Extracting months and weekdays from data:

```
> head(months(fd$FlightDate))
[1] "January" "January" "January" "January" "January"
> unique(months(fd$FlightDate))
[1] "January" "February" "March"
> head(weekdays(fd$FlightDate))
[1] "Wednesday" "Wednesday" "Wednesday" "Wednesday" "Wednesday"
> unique(weekdays(fd$FlightDate))
[1] "Wednesday" "Thursday" "Friday" "Saturday" "Sunday" "Monday" "Tuesday"
```

Computing time intervals and using difftime()

```
> fd$FlightDate[60]-fd$FlightDate[900]
Time difference of -3 days
> difftime(fd$FlightDate[3000],fd$FlightDate[90],units = "weeks")
Time difference of 1.571429 weeks
> difftime(fd$FlightDate[3000],fd$FlightDate[90],units = "days")
Time difference of 11 days
> difftime(fd$FlightDate[3000],fd$FlightDate[90],units = "hours")
Time difference of 264 hours
```

1

- Manipulating data involving dates
- Sub-setting data: All rows when the day is Sunday

```
> library(dplyr)
> #Subset the data for day=Sunday
> dim(fd)
[1] 30443     25
> fd_s<-fd%>%filter(weekdays(FlightDate)=="Sunday")
> dim(fd_s)
[1] 4015     25
```

- Manipulating data involving dates
- Find the number of flights on Sundays for destination Atlanta

```
> #Find the number of flights on Sundays for destination Atlanta
```

> fd%>%filter(weekdays(FlightDate)=="Sunday",DestCityName=="Atlanta, GA")%>%nrow()
[1] 683

-

- Manipulating data involving dates
- Find the number of flights on Sundays for all cities

```
> #Find the number of flights on Sundays for all cities
> fd%>%filter(weekdays(FlightDate)=="Sunday")%>%group_by(DestCityName)%>%summarize(n())
Source: local data frame [10 x 2]

        DestCityName n()
1        Atlanta, GA 683
2        Charlotte, NC 342
3        Chicago, IL 193
4        Denver, CO 448
5        Houston, TX 155
6        Las Vegas, NV 507
7        Los Angeles, CA 603
8        New York, NY 349
9        Phoenix, AZ 466
10 Washington, DC 269
```

 Whenever data has time information along with date, R uses POSIXct and POSIXIt classes to deal with dates

```
> date1<-Sys.time()
> date1
[1] "2015-03-02 17:35:47 IST"
> class(date1)
[1] "POSIXCT" "POSIXT"
> weekdays(date1)
[1] "Monday"
> months(date1)
[1] "March"
```

 Whenever data has time information along with date, we use POSIXct and POSIXIt classes to deal with dates

```
> date2<-as.POSIXlt(date1)</pre>
> date2
[1] "2015-03-02 17:35:47 IST"
> str(date2)
 POSIXIt[1:1], format: "2015-03-02 17:35:47"
> date2$wday
[1] 1
> date2$zone
[1] "IST"
> date2$hour
[1] 17
> date2$wday
> date2$zone
[1] "IST"
> date2$hour
[1] 17
```

- lubridate() is a package that is a wrapper for POSIXct class
- It has a very simple syntax

```
> library(lubridate)
> fd$FlightDate<-dmy(fd$FlightDate)</pre>
> str(fd)
'data.frame': 30443 obs. of 25 variables:
 $ FlightDate
                     : POSIXct. format: "2014-01-01" "2014-01-01" "2014-01-01" ...
 $ UniqueCarrier
                     : Factor w/ 13 levels "AA", "AS", "B6",...: 4 4 4 4 11 11 11 11 11 4 ...
 $ AirlineID
                     : int 19790 19790 19790 19790 20355 20355 20355 20355 19790 ...
                     : Factor w/ 13 levels "AA", "AS", "B6",..: 4 4 4 4 11 11 11 11 11 4 ...
 $ Carrier
                     : Factor w/ 2816 levels "", "D942DN", "N0EGMQ",...: 2641 2512 2490 2581 1657 34 60 17
 $ TailNum
98 1443 473 ...
 $ FlightNum
                     : int 335 1095 2422 1607 657 894 1843 2041 413 2030 ...
 $ OriginAirportID
                     : int 11057 11057 11057 11057 11057 11057 11057 11057 11057 13232 ...
 $ OriginAirportSeqID: int 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 1105703 132
3202 ...
 $ originCityMarketID: int 31057 31057 31057 31057 31057 31057 31057 31057 31057 31057 31057 31057
                     : Factor w/ 10 levels "ATL", "CLT", "DCA", ...: 2 2 2 2 2 2 2 2 2 9 ...
 $ Origin
 $ OriginCityName
                     : Factor w/ 10 levels "Atlanta, GA",...: 2 2 2 2 2 2 2 2 3 ...
                     : Factor w/ 10 levels "AZ", "CA", "CO", ...: 6 6 6 6 6 6 6 6 5 ...
 $ OriginState
```

- lubridate() is a package that is a wrapper for POSIXct class
- It has a very simple syntax.

Function	Date
dmy()	26/11/2008
ymd()	2008/11/26
mdy()	11/26/2008
dmy_hm()	26/11/2008 20:15
dmy_hms()	26/11/2008 20:15:30

RECAP

- Character to date conversion-Date Class
- Extracting months and weekdays
- Using difftime()
- Manipulating data involving dates
- POSIXct and POSIXIt Classes
- Working with lubridate()

Class 4 – Data Manipulation in R

Topic 4



Merging Tables



INDEX

Manipulating data using base R Using dplyr to manipulate data Working with date objects



Merging tables

Missing value treatment
Using reshape2() to transpose data
Manipulating Character Strings
Using sqldf





Joining dataframes

- Just like tables can be joined in sql, we can perform joins on dataframes in R
- Following types of joins can be accomplished
- Inner Join
- Left outer join
- Right outer join
- Full outer join



Joining dataframes

Inner join: Joining two tables based on a key column, such that rows matching in both tables are selected

	CustomerId	Product		CustomerId	State
1	1	Toaster	1	2	Alabama
2	2	Toaster	2	4	Alabama
3	3	Toaster	3	6	Ohio
4	4	Radio	_	-	
5	5	Radio			
6	6	Radio			

-

Joining dataframes

 Inner join: Joining two tables based on a key column, such that rows matching in both tables are selected

	CustomerId	Product
_1	1	Toaster
2	2	Toaster
3	3	Toaster
4	4	Radio
5	5	Radio
6	6	Radio

```
CustomerId State
1 2 Alabama
2 4 Alabama
3 6 Ohio
```

```
> merge(x=df1,y=df2,by="CustomerId")#Inner Join/Intersection of both tables
  CustomerId Product State
1     2 Toaster Alabama
2     4 Radio Alabama
3     6 Radio Ohio
```

Joining dataframes

 Full Outer Join: Two tables are joined irrespective of any match between the rows

```
CustomerId Product
                             CustomerId
                                          State
           1 Toaster
                                      2 Alabama
           2 Toaster
                                      4 Alabama
           3 Toaster
                                           Ohio
           4 Radio
           5 Radio
           6 Radio
> merge(x = df1, y = df2, by = "CustomerId", all = TRUE)#Outer join:
 CustomerId Product
                     State
          1 Toaster
                      < NA >
          2 Toaster Alabama
                    <NA>
          3 Toaster
          4 Radio Alabama
          5 Radio <NA>
          6 Radio Ohio
```

1

Joining dataframes

 Left Outer Join: All the rows of left table are retained while matching rows of right table are displayed

```
CustomerId Product
                              CustomerId
                                            State
           1 Toaster
                                       2 Alabama
           2 Toaster
                                       4 Alabama
           3 Toaster
                                             Ohio
           4 Radio
           5 Radio
           6 Radio
> merge(x = df1, y = df2, by = "CustomerId", all.x=TRUE)#Left join
 CustomerId Product
                      State
          1 Toaster
                       < NA >
          2 Toaster Alabama
          3 Toaster
                      <NA>
          4 Radio Alabama
          5 Radio
                       <NA>
          6 Radio
                    Ohio
```

Joining dataframes

 Right Outer Join: All the rows of right table are retained while matching rows of left table are displayed

```
CustomerId Product CustomerId State

1 1 Toaster 1 2 Alabama
2 Toaster 2 4 Alabama
3 Toaster 3 6 Ohio
4 Radio
5 Radio
6 Radio
```

```
> merge(x = df1, y = df2, by = "CustomerId", all.y=TRUE)#Right join
CustomerId Product State
1    2 Toaster Alabama
2    4 Radio Alabama
3    6 Radio Ohio
```

RECAP

- Inner Join
- Left outer join
- Right outer join
- Full outer join



Topic 5



Missing Value Treatment



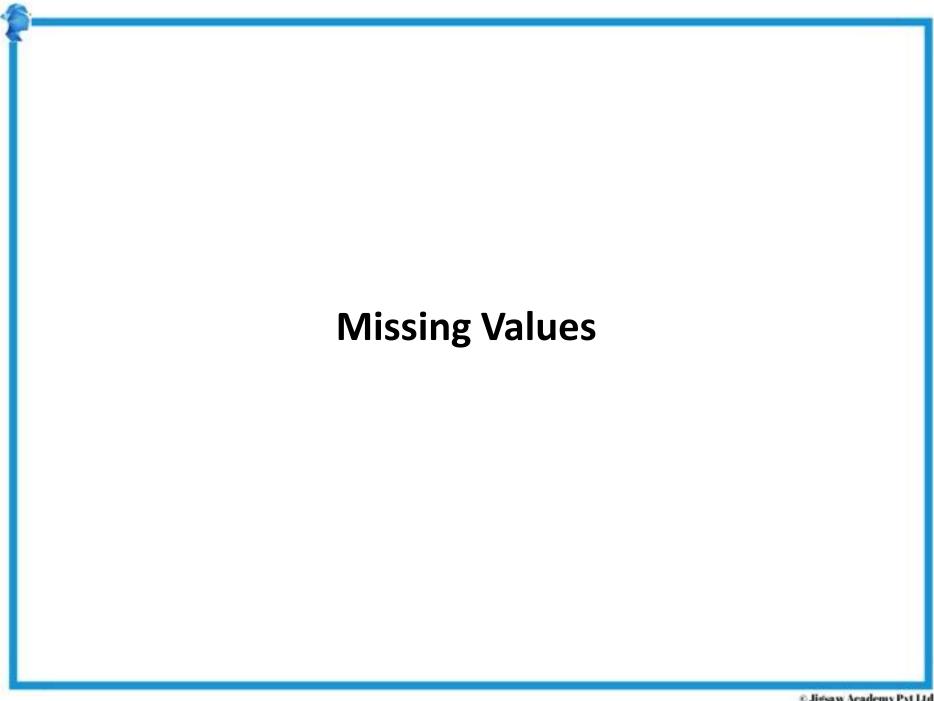
INDEX

Manipulating data using base R
Using dplyr to manipulate data
Working with date objects
Merging tables



Missing value treatment

Using reshape2() to transpose data Manipulating Character Strings Using sqldf





- Identifying missing values in a column
- Imputing missing values



Using is.na() to find out the total number of missing values

```
> a<-c(1,2,3,4,5,6,NA,NA,NA,7,8,9)
> is.na(a)
[1] FALSE FALSE FALSE FALSE FALSE TRUE TRUE TRUE FALSE FALSE FALSE
> sum(is.na(a))
[1] 3
```

-

Missing Values

Using is.na() to find out the total number of missing values

```
> air<-airquality
> head(air)
 Ozone Solar.R Wind Temp Month Day
     41
           190
                       67
     36
           118 8.0
                      72
    12
           149 12.6
                      74
4
5
                      62 5
    18
           313 11.5
            NA 14.3
                       56
     NA
     28
            NA 14.9
                      66
```

Using is.na() to find out the total number of missing values

```
> sum(is.na(air$0zone))
[1] 37
> sum(is.na(air$Solar.R))
[1] 7
> summary(air)
                     Solar.R
                                       Wind
                                                                        Month
     Ozone
                                                         Temp
                                                                                          Day
                                                           :56.00
Min.
      : 1.00
                  Min. : 7.0
                                         : 1.700
                                                    Min.
                                                                    Min.
                                                                           :5.000
                                                                                    Min.
                                  Min.
                                                                                          : 1.0
1st Qu.: 18.00
                                  1st Qu.: 7.400
                  1st Qu.:115.8
                                                    1st Qu.:72.00
                                                                    1st Qu.:6.000
                                                                                    1st Qu.: 8.0
Median : 31.50
                                  Median : 9.700
                                                    Median:79.00
                  Median :205.0
                                                                    Median :7.000
                                                                                    Median:16.0
Mean : 42.13
                         :185.9
                                        : 9.958
                                                           :77.88
                                                                           :6.993
                                                                                          :15.8
                  Mean
                                  Mean
                                                    Mean
                                                                    Mean
                                                                                    Mean
 3rd Qu.: 63.25
                  3rd Qu.:258.8
                                  3rd Qu.:11.500
                                                    3rd Qu.:85.00
                                                                    3rd Qu.:8.000
                                                                                    3rd Qu.:23.0
        :168.00
                         :334.0
                                          :20.700
                                                           :97.00
                                                                           :9.000
                                                                                            :31.0
Max.
                  Max.
                                  Max.
                                                    Max.
                                                                    Max.
                                                                                    Max.
NA's
        :37
                  NA's
                         :7
```



Imputing missing values

- > #Imputing Missing values
 > air\$ozone[is.na(air\$ozone)]<-45
 > summary(air)

2.5	_				
Ozone	Solar.R	Wind	Temp	Month	Day
Min. : 1.00	Min. : 7.0	Min. : 1.700	Min. :56.00	Min. :5.000	Min. : 1.0
1st Qu.: 21.00	1st Qu.:115.8	1st Qu.: 7.400	1st Qu.:72.00	1st Qu.:6.000	1st Qu.: 8.0
Median : 45.00	Median :205.0	Median : 9.700	Median :79.00	Median :7.000	Median :16.0
Mean : 42.82	Mean :185.9	Mean : 9.958	Mean :77.88	Mean :6.993	Mean :15.8
3rd Qu.: 46.00	3rd Qu.:258.8	3rd Qu.:11.500	3rd Qu.:85.00	3rd Qu.:8.000	3rd Qu.:23.0
Max. :168.00	Max. :334.0	Max. :20.700	Max. :97.00	Max. :9.000	Max. :31.0
	NA'S :7				



Imputing missing values

- > #Imputing Missing values
 > air\$Solar.R[is.na(air\$Solar.R)]<-mean(air\$Solar.R,na.rm=TRUE)</pre>
- > summarv(air)

Ozone	Solar.R	Wind	Temp	Month	Day
Min. : 1.00	Min. : 7.0	Min. : 1.700	Min. :56.00	Min. :5.000	Min. : 1.0
1st Qu.: 21.00	1st Qu.:120.0	1st Qu.: 7.400	1st Qu.:72.00	1st Qu.:6.000	1st Qu.: 8.0
Median : 45.00	Median :194.0	Median : 9.700	Median :79.00	Median :7.000	Median :16.0
Mean : 42.82	Mean :185.9	Mean : 9.958	Mean :77.88	Mean :6.993	Mean :15.8
3rd Qu.: 46.00	3rd Qu.:256.0	3rd Qu.:11.500	3rd Qu.:85.00	3rd Qu.:8.000	3rd Qu.:23.0
Max. :168.00	Max. :334.0	Max. :20.700	Max. :97.00	Max. :9.000	Max. :31.0

RECAP

- Identifying missing values in a column
- Imputing missing values

Class 4 – Data Manipulation in R

Topic 6

**ansposing data using reshape2

INDEX

Manipulating data using base R

Using dplyr to manipulate data

Working with date objects

Merging tables

Missing value treatment



Using reshape2() to transpose data

Manipulating Character Strings

Using sqldf





- Understanding wide and long data formats
- Converting data in wide format to long
- Converting data in long format to wide



Transposing data

- Understanding wide and long data formats
- Most structured data is in wide format: Variables are columns and Row labels identify observations

Person	Age	Weight
Sankar	26	70
Aiyar	24	60
Singh	25	65

Transposing data

- Understanding wide and long data formats
- The same data can be represented as follows (long format)

Persons	Variable	Value
Sankar	Age	26
Sankar	Weight	70
Aiyar	Age	24
Aiyar	Weight	60
Singh	Age	25
Singh	Weight	65

-

Transposing data

Converting a wide format data to long format: melt()

```
> library(reshape2)
> person<-c("Sankar","Aiyar","Singh")
> age<-c(26,24,25)
> weight<-c(70,60,65)
> wide<-data.frame(person,age,weight)
> wide
   person age weight
1 Sankar 26 70
2 Aiyar 24 60
3 Singh 25 65
```



Transposing data

Converting a wide format data to long format: melt()

```
> melt(wide)
Using person as id variables
  person variable value
1 Sankar
              age
                     26
2 Aiyar
                     24
              age
3 Singh
                     25
              age
4 Sankar
           weight
                     70
5 Aiyar
           weight
                     60
6 Singh
           weight
                     65
> melt(wide,id.vars="person")
  person variable value
1 Sankar
                     26
              age
2 Aiyar
                     24
              age
3 Singh
              age
                     25
4 Sankar
           weight
                     70
5 Aiyar
           weight
                     60
6 Singh
           weight
                     65
```

-

Transposing data

Converting a wide format data to long format: melt()

```
> melt(wide,id.vars="person",variable.names="Demographics",value.name ="Demo_Value" )
  person variable Demo_Value
1 Sankar
                          26
              age
2 Aiyar
             age
3 Singh
              age
4 Sankar
          weight
5 Aiyar
          weight
6 Singh
         weight
                         65
> melted<-melt(wide,id.vars="person",variable.names="Demographics",value.name ="Demo_value")
```



Transposing data

Converting a long format data to wide format: dcast()

```
> dcast(melted,person~variable,value.var = "Demo_Value")
  person age weight
1 Aiyar 24 60
2 Sankar 26 70
3 Singh 25 65
```

RECAP

- Understanding wide and long data formats
- Converting data in wide format to long
- Converting data in long format to wide



Topic 7



Manipulating Character Strings and ** Using sqldf



INDEX

Manipulating data using base R

Using dplyr to manipulate data

Working with date objects

Merging tables

Missing value treatment

Using reshape2() to transpose data



Manipulating Character Strings

Using sqldf



-

Manipulating Strings in R

- We'll cover the following R functions used to manipulate character data
- substr()
- nchar()
- tolower(), toupper()
- strsplit(),paste()
- grep(),grepl(),
- sub(),gsub()



Manipulating Strings in R

substr(), nchar(),tolower(),toupper()

```
> a<-"Batman"
> substr(a,start=2,stop=6)
[1] "atman"
> nchar(a)
[1] 6
> tolower(a)
[1] "batman"
> toupper(a)
[1] "BATMAN"
```

-

Manipulating Strings in R

strsplit(), paste(),grep(),grepl(),sub(),gsub()

```
> b<-"Bat-Man"
> c<-"Bat/Man"
> strsplit(c,split="/")
[[1]]
[1] "Bat" "Man"
> paste(b,split=c)
[1] "Bat-Man Bat/Man"
> grep("-",b)
[1] 1
> grep1("/",c)
[1] TRUE
> sub("-","/",b)
[1] "Bat/Man"
> d<-"Bat-Ma-n"
> sub("-","/",d)
[1] "Bat/Ma-n"
> gsub("-","/",d)
[1] "Bat/Ma/n"
```

INDEX

Manipulating data using base R

Using dplyr to manipulate data

Working with date objects

Merging tables

Missing value treatment

Using reshape2() to transpose data

Manipulating Character Strings



Using sqldf



SQL queries within R: sqldf

- One can easily use sql queries within R using a package called sqldf
- This package is available at CRAN
- install.package(sqldf)
- library(sqldf)

-

Manipulating Strings in R

Using select statement

```
> library(sqldf)
> #Using SELECT statement
> oj_s<-sqldf("select brand, income, feat from oj ")</pre>
```

- Using where statement
- > #Subseting using where statement
 > oi s<-sqldf("select brand. income. feat from oi where price<3.8 and income<10")</pre>
- Using order by statement

```
> #Order by statement
> oj_s<-sqldf("select store,brand,week,logmove,feat,price, income from oj order by income asc")</pre>
```



SQL queries within R: sqldf

Using sql functions

```
> #Demo sql functions
> sqldf("select avg(income) from oj")
  avg(income)
1    10.61673
> sqldf("select min(price) from oj")
  min(price)
1    0.52
```

RECAP

- Basic string manipulation functions: tolower, toupper,grep, grepl, sub, gsub
- Executing sql queries in R, using sqldf() package