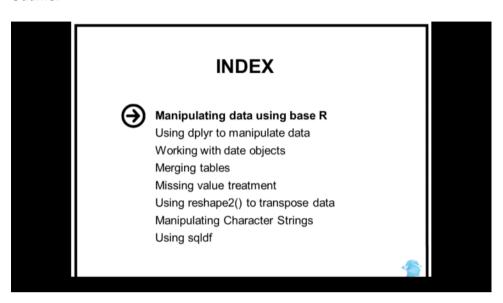
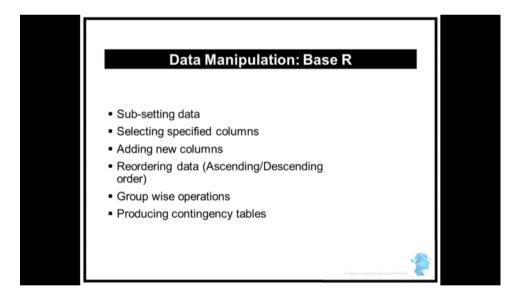
## Manipulating Data with R

Tuesday, August 23, 2016 4:41 PM

### Outline:





To access one specific cells in a data frame we use the following command --> df[1,3] To access many specific cells in a data frame we can use the following command --> df[c(1,4,6,198),c(4,5,9)]

We can also access using a specific column name --> df[c(1:5),"brand"]

Lets look at logical subsetting of data i..e

If we want to get a subset of data, using a column filter:

Find all rows of data where band column has value as "Tropicana" --> df[df\$Brand == 'Tropicana',]

Hence the condition is in the rows part of the DF in the code.

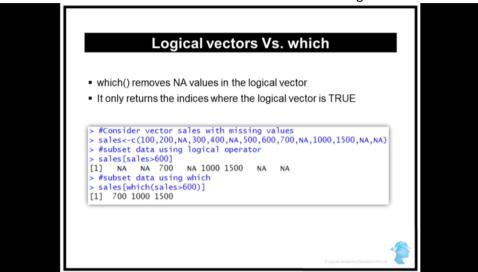
If we want to use multiple filters , we can do that by using the OR operator of R i.e. | Df[df\$brand=='Tropicana'|df\$brand=='minutemaid',]

We can also combine data output based on AND operators ---> i.e. & Df[df\$brand=='Tropicana' & df\$feat==0,]

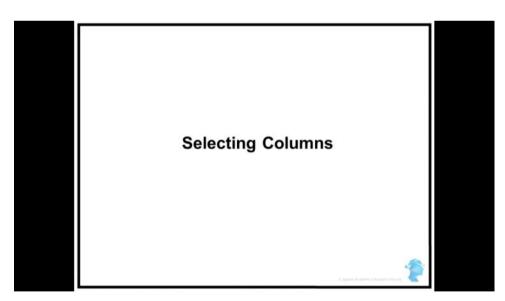
We can use these combinations to subset the data.

Another way to subbset the data is using the which statement: this will give row indices Index<-which(df\$brand=='Tropicana')

What is the difference between a Which statement and a logical statement:



Logical statements will include the "NA" values but which will only return the values that are true for the condition



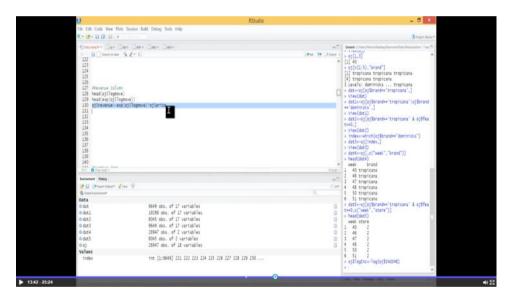
To select specific columns, we can mention either the column index or column name in our query : Df[, c("brand", "feat")]

To subset based on logical conditions we can do the following: Df[df\$brand=="tropicana" & df\$feat==0, c("week","store")]

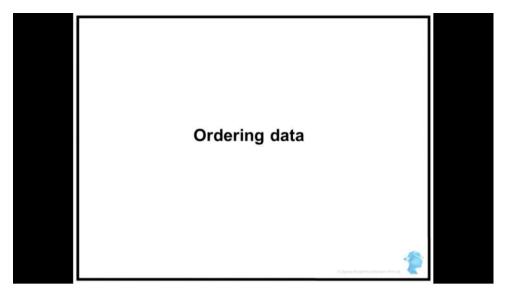


To add a new column to an existing data frame we can use this: Df\$logInc<-log(df\$income)

We can store logical operations in new columns as follows: Df\$income<-exp(df\$logmove)\*df\$price



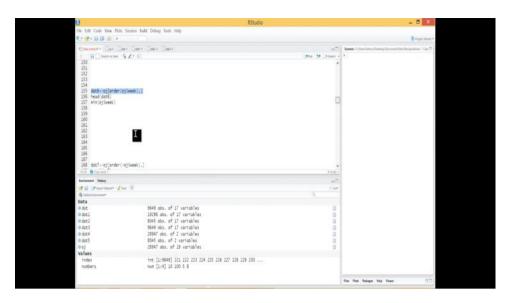
New is , how to sort the data:

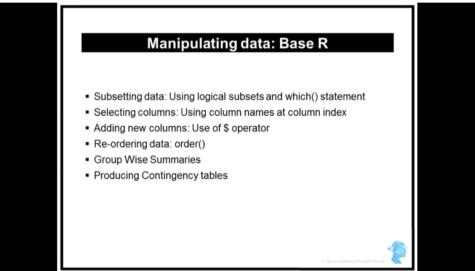


## Ordering order() returns the element order that results in a sorted vector students students [1] "John" "Tim" "Alice" "Zeus" order(students) [1] 3 1 2 4 students[order(students)] [1] "Alice" "John" "Tim" "Zeus" Application: Very useful for sorting dataframes

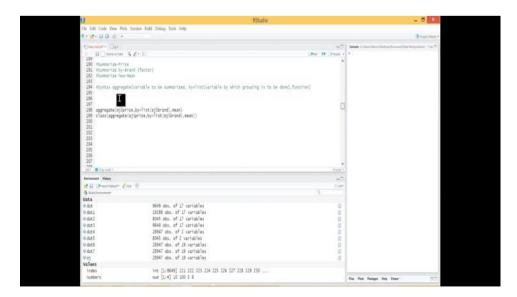
We can use the base R function called: order
Just using order(students) --> this will give the indices
Using --> students[order(students)] --> this will give the dataset
If we want to sort our data frame based on week numbers:
df[order(df\$week),]

To reverse order the dataframe --> df[order(-df\$week),]





The next function is group by summary: unique Df[(unique(df\$brand)]



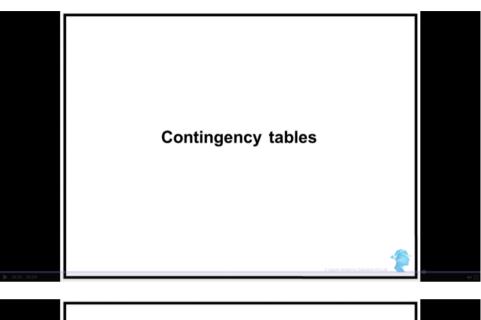
Syntax for aggregation:

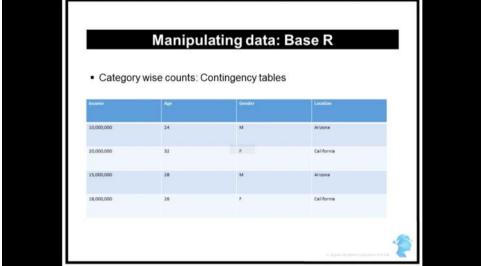
Aggregate(variable to be summarized, by=list(variables by which grouping is to be done), function)

Aggregate(df\$price, by=list(df\$brand), mean)

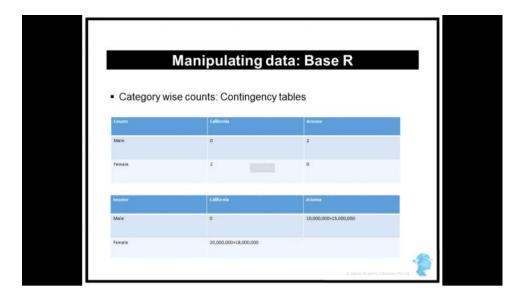
This is always of the class data frame

We can perform the same functionality by using tapply Tapply(df\$price, df\$brand, mean)
In this case the class of the object is an array

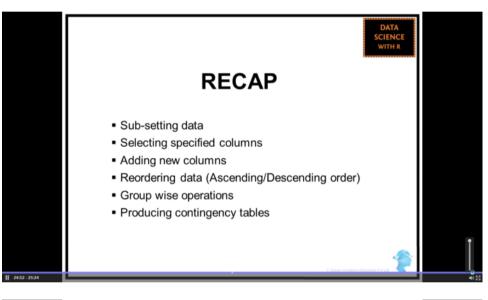




Contingency tables are those that have categorical values



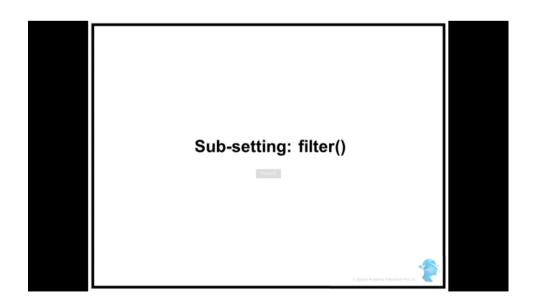
Cross tabulation , use the table command --> table(df\$brand, df\$feat)



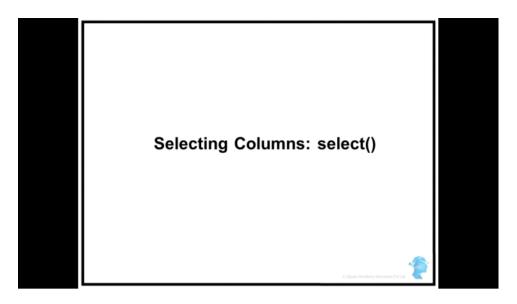


# Manipulating data: dplyr dplyr: Whats and Whys Sub-setting data using filter() Selecting columns using select() Adding new columns using mutate() Ordering data using arrange() Summarizing using summarize() and group\_by() Using functional pipelines to do more than one manipulation task

## 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)

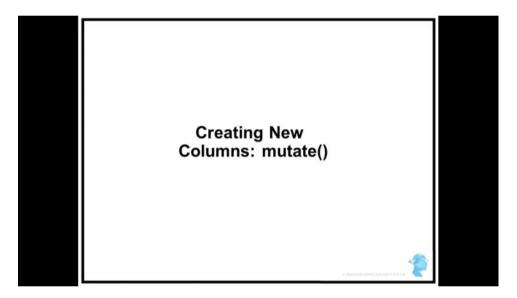


Data<-filter(df, brand == "tropicana"|brand=="minutemaid")</pre>

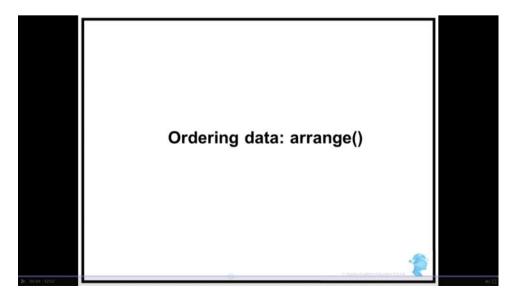


data<-select(df, brand, income, feat)

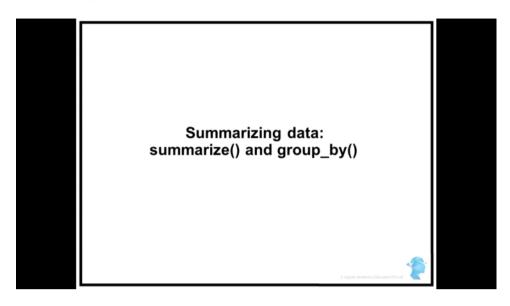
We can use the -ve mark if we want to exclude culumns from the original set into the new dataset Data<-select(df, -brand, -feat)



Data<-mutate(df, logincome=log(income))

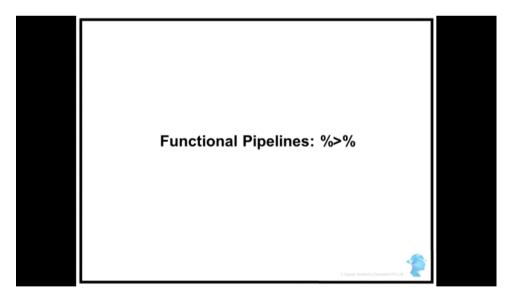


Data<-arrange(df,income)
Data<-arrange(df,desc(income))

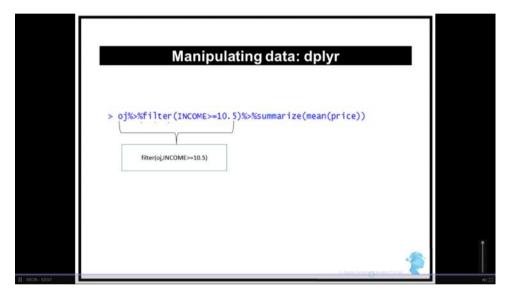


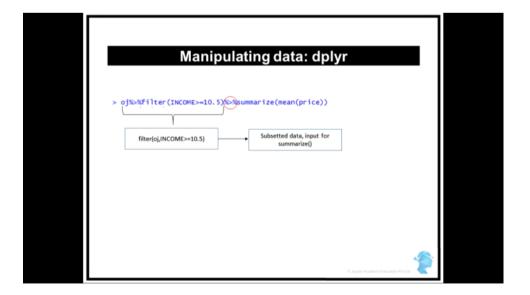
Gr\_brnad<-group\_by(df,brand)
Summarize<-(gr\_band, mean(income), sd(income))</pre>

For the summarize function , the first attribute is the input and the rest are the outputs(function that you want to carry on the input).

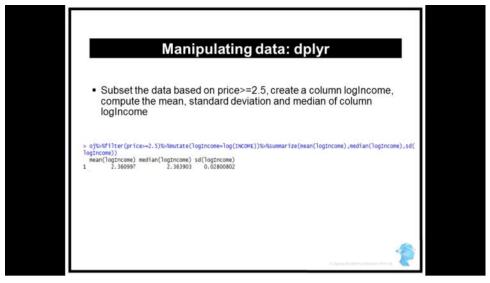


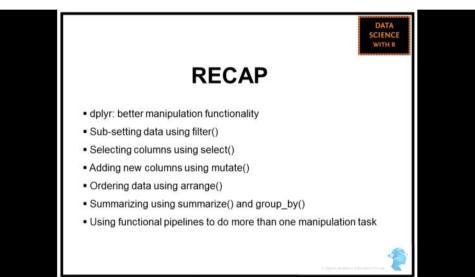
```
    Manipulating data: dplyr
    dplyr becomes a powerful tool when combined with %>% (pipe) operator
    Several data manipulation tasks can be accomplished in just one line of code
    Traditionally functional composition is achieved by using nested function calls
    For example, Find the mean price for all people whose income is >=10.5
    #Base R code > mean(oj[ojSINCOME>=10.5, "price"]) [1] 2.770229
    #dplyr code > summarize(filter(oj,INCOME>=10.5), mean(price)) mean(price) 1 2.270229
```



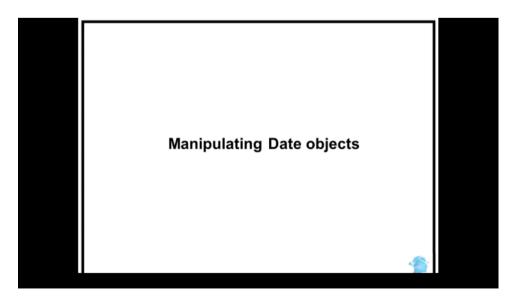


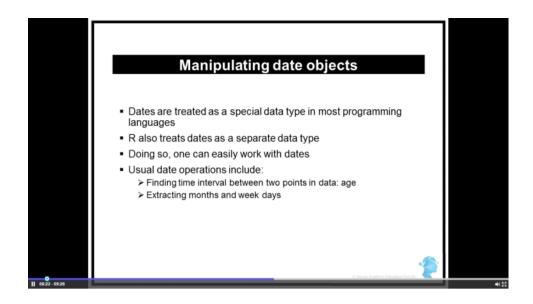
<Make sure that the data type is compatible while using pipes i.e. %>%

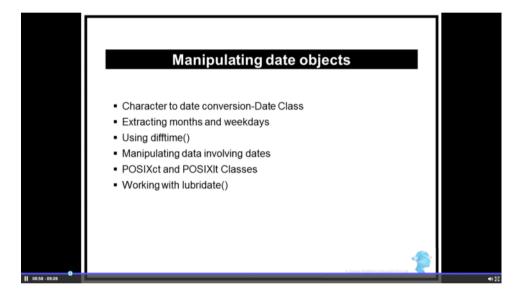


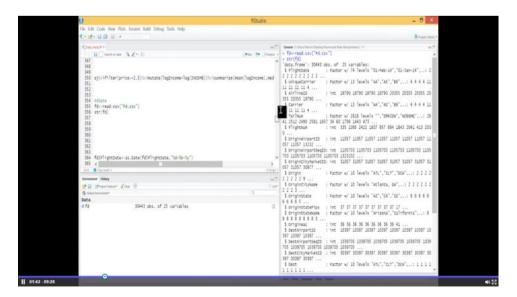


## Working with Date Objects:

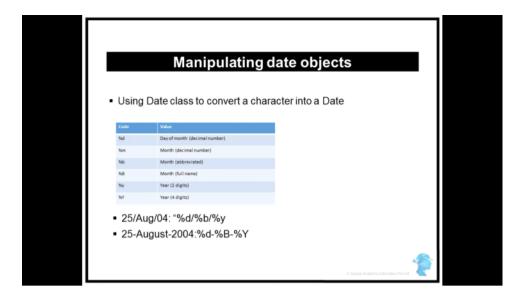






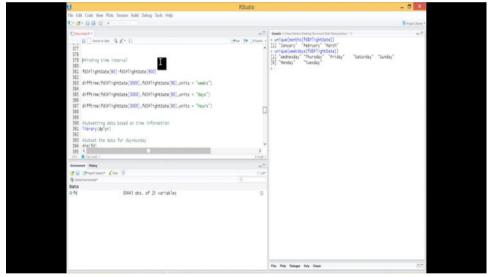


Here are the different data formats for date:



To extract the month from the dates use --> months(fd\$flightdate)
To extract unique months --> unique(months(fd\$flightdate))
Similarly to get weekdays --> unique(weekdays(fd\$flightdate))

To find the time interval between two dates:

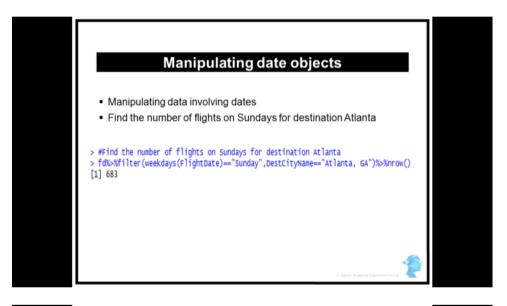


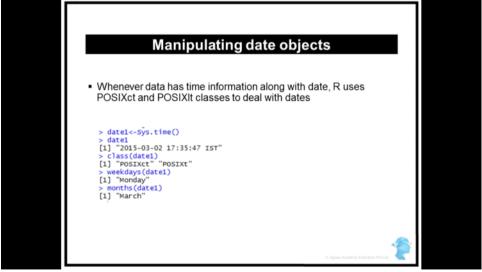
The difftime function can be used to find the timedifferences , the units can be in days, minutes, weeks etc

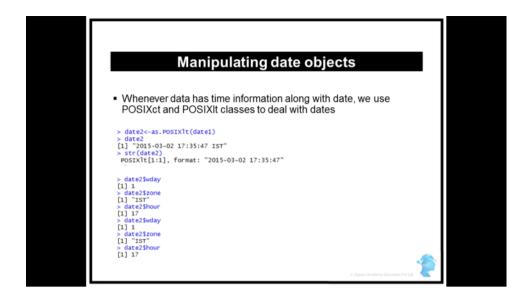
Difftime(fd\$Flightdate[3000],fd\$flightdate[60], units="weeks")

# Manipulating date objects • 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<-fdb>%filter(weekdays(FlightDate)=="Sunday") > dim(fd\_s) [1] 4015 25

Using DPLYR for the similar functionalities with dates:



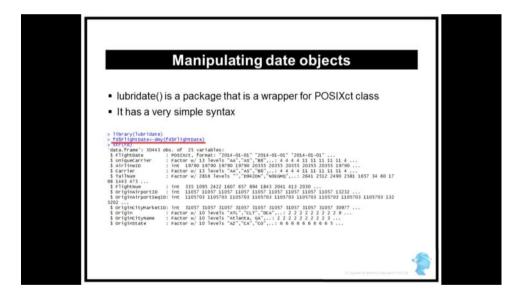


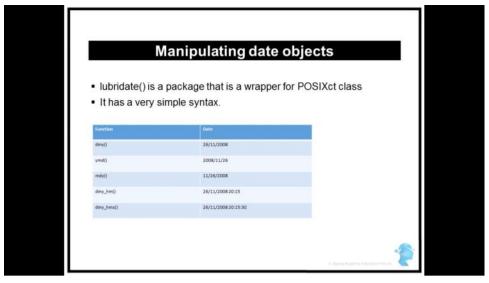


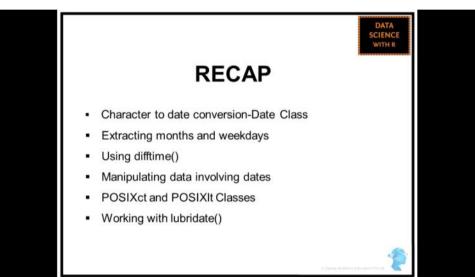
### Date2<-as.POSIXIt(date1)

This will store the various date attributes stored as list and makes it easier to exctract this data

Date2\$wday Date2\$zone Date2\$hour







## Merging Tables:





### 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

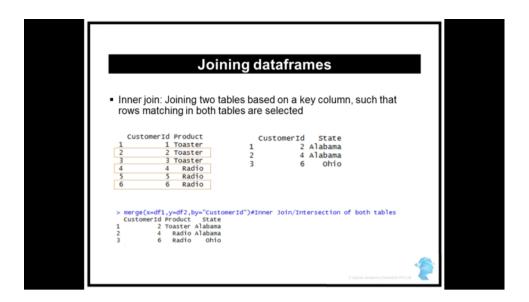


### Joining dataframes

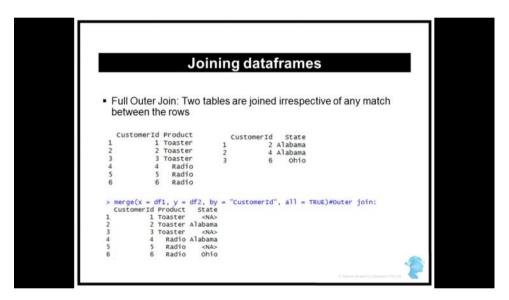
 Inner join: Joining two tables based on a key column, such that rows matching in both tables are selected



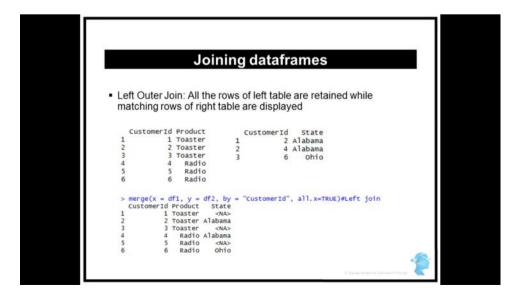




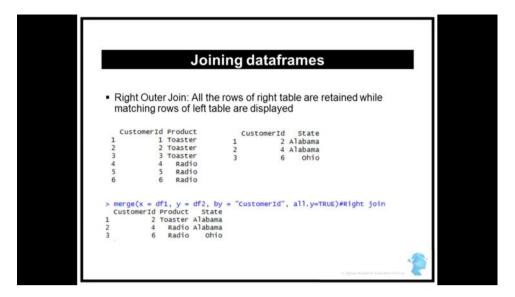
To perform inner join we use the following <- merge(x=df1, y=df2, by="customer id")



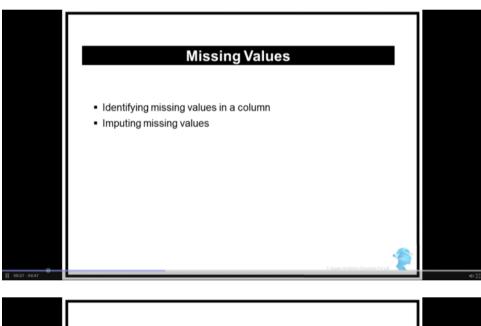
For full outer join <- merge(x=df1, y=df2, by="customer id", all= TRUE)

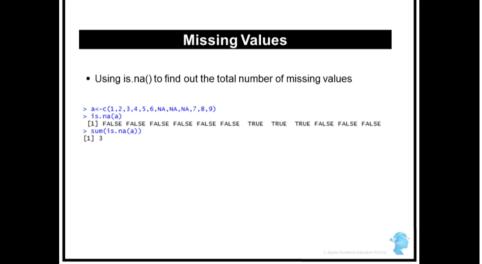


Left outer join <- merge(x=df1, y=df2, by="customer id", all.x= TRUE) right outer join <- merge(x=df1, y=df2, by="customer id", all.y= TRUE)

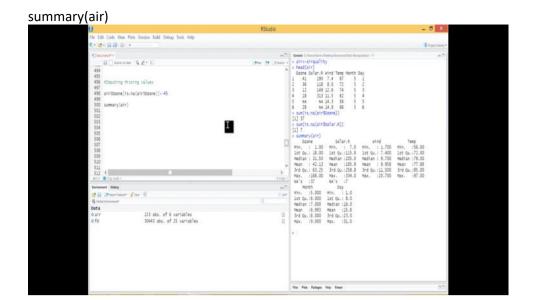


## **Treating Missing Values:**



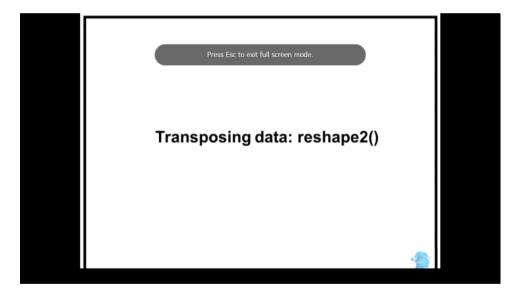


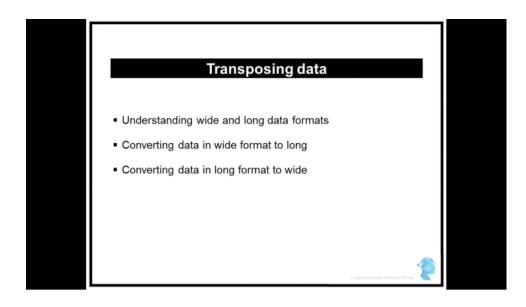
Running a summary command will also give the number of missing values in each column :-

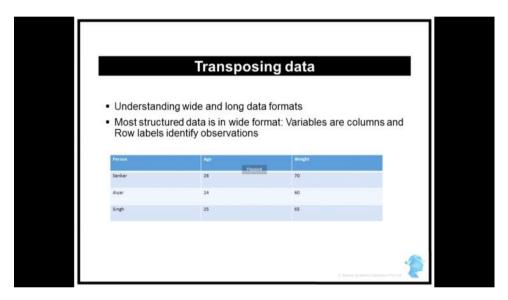


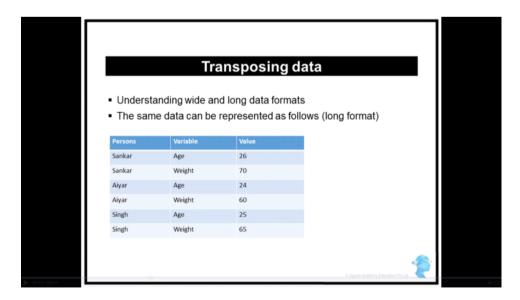
Imputing these missing values :
Air\$ozone[is.na(air\$ozone)]<-45
Air\$Solar.R[is.na(air\$Solar.R)]<-mean(air\$Solar.R,na.rm=TRUE)

## Transposing Data using Reshape2:









Use the library --> reshape2

Create a data in wide or long format

Wide <- data.frame(person, age, weight)

To convert this to long format we will use the function called Melt from reshape2

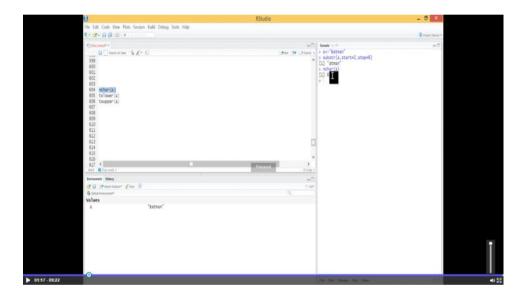
Melt(wide,id.vars="person", value.name="demo\_value") This will generate the long format.

To convert data from long format to wide format we will use dcast

Dcast(melted, person~variable, value.var="demo\_value")

## Manipulating character strings and using SQLDF

String amnipulation, we can use substr a<-"Batman"
Susbstr(a,start=2, stop=6)



To find the number of characters in a string --> nchar(a)
To convert all to lower case --> tolower(a)
TO covert all to upper case --> toupper(a)

To split a string, mention the delimiter to split with Strsplit(a, split="-")

To concatenate two objects together we use paste Paste(b, split=c)

If we want to find a character in the elements of a vector we can use : grep Grep("-", c(a,b))

We can also use: grepl

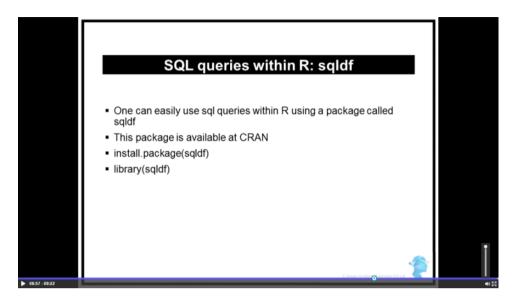
This does a logical comparision and returns the results for all the elements based on satisfying the conndition.

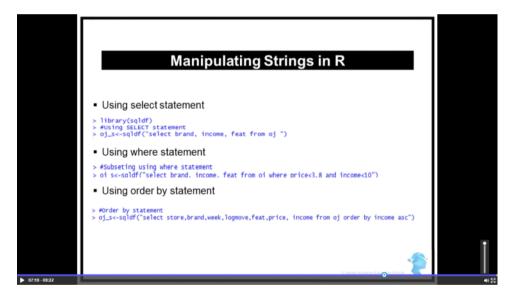
Grepl("-",c(a,b))

If we want to substitute a character with another character in an element we use: sub

Sub("-","/",a)

But this command only replaces the first occurance of the string. TO replace all the occurances: Use gsub("-","/",a)





Using the select statement is very easy with this SQLDF library:

Obj\_sql<-sqldf("select \* from a")

Make sure the object is a table to use with this select statement.

