Will Speak too fast and mumble words

Please ask Siva to repeat if you do not understand

• R is an open-source compute programming language and runs on Linux, Windows, and Mac-OS

• R is FREE.

• The R project web page http://www.r-project.org

on Lin	ux, Windows, and Mac-OS.	
\$ -	www - help	
	araj lable	
•	R - bbggca	

• R is modeled after S and S-Plus. The S language was developed in the late 1980s at AT & T labs.

• The R project was started by Robert Gentleman and Ross Ihaka of the Statistics Department of the University of Auckland in 1995. [Journal of Computational and Graphical Statistics,5:3, pp. 299-314. 1996.]

• R is now a collaborative project with many contributors and is maintained by the R core-development team.

• To download R visit https://cloud.r-project.org

• Rstudio is an Integrated Development Environment, or IDE for R. To download Rstudio visit http://www.rstudio.com/download

	$\int 0 \ 0 \ 1$	
	CHOND	
 Data and its analysis has a rich and wide literature. Three kinds of Data: 		
 Categorical Data Discrete Numeric Data Continuous Numeric Data 	Numbers!	
• On the Theory of Scales of Measurement By S. S. Ste gave a broad classification of data from measurement	evens Science 07 Jun 1946: Vol. 103, Issue 2684, pp. 677-680, s into 9 categories.	

• Many data are described in terms of numbers.

• Many variables naturally take on only discrete values.

• Boxplot and Histograms are used to visualise such data.

Daily data of Temperature
Earnings in a stock-market

Discrete Numerical Data: Key features



- Center
- Spread
- Shape

Discrete Numerical Data: Key features



- <u>Center</u> Widely used measure of centre is the <u>mean</u> or the average of the data set. Other measures include the median and the mode
- <u>Spread</u> Understanding variability of the given data is very important. If one were to understand mean as specifying the center then the range of the data set around it is determined by its variability or spread. It is often measured by the variance(var) or standard deviation (sd) or the inter-quartile range (IQR).
- Shape To understand various distributional aspects of the dataset one needs to understand its "shape". For e.g. if it is symmetric or skewed around its mean. Other aspects include among the data points which are more likely than others.

very low chance of occuring Range Essential × K nean + 3.5.0 mean -3SD megn

- Ubeful collection - Strongly encourage to find a data set R has a lot inbuilt Datasets that one can use. The command : > data() will list currently installed data sets. in a class. R has a lot inbuilt Datasets that one can use. The command :
> data()

will list currently installed data sets.

Pords	2	Compre hensiv	C	
list.		•		

- R stores many datasets as data frame (often).
- A data frame is a rectangular collection of variables (in the columns) and observations (in the rows).

Let us learn about real data stored as data frame.

> ?airquality

Describes ltic data set	
air quality	
· 153 - Rows	
6 - Columns	

Let us learn about airquality dataset a bit more.

- we could print the entire data set on the screen
 >airquality
 but this is too much information.
- Let us try the head() function
 - > head(airquality)

This provides the first six rows.

Let us learn about airquality dataset a bit more.

- Let us try the tail() function
 - > tail(airquality)

This provides the last six rows.

• Below provides the first ten rows.	
> head(airquality, $n = 10$)	
Data can be called using row and column number	
> airquality[148,4]	
[1] 63	
• We can use the variable name for the given column and call it by its position.	
> airquality\$Temp[148]	
[1] 63	

Provides an entire row
> airquality[148,]
Provides Ozone Temp columns
> airquality[,c(1,4)]
using c() function we can form any vector and that will enable display of the respective columns. We did not specify the row, so all rows will be displayed.

Five Number Summary and Histograms





Plot

We can use the plot function to just plot.

> plot(airquality\$Temp)



Scatter Plot

We can use the plot function to get a Scatter plot.
> plot(airquality\$0zone, airquality\$Temp)



> plot(airquality)



R has can be enhanced with a lot of external packages that are available. The package UsingR has many datasets loaded in it.

```
> install.packages("UsingR")
```

Once installed then to add to current workspace

> library("UsingR")

ggplot2	implements	grammar	of	graphics
66012	implements	grannar		graphics

> install.packages("tidyverse")

Once installed then to add to current workspace

> library("tidyverse")

gramma of graphics	
- system for describing	
2 building graphs	
- learn one system and	
apply it at many	
places.	

Dataset in tidyverse		
> mpg		
Observations collected by US Environment Protection Agency on 38	models of cars.	

ggplot2- Data Visualisation



• Begins with a function ggplot()- creates a coordinate system that you can add addlayers to. The first arugment is the data set to use ggplot(data=mpg) creates an empty graph. • Add layers to ggplot()- the function geom_point() adds a layer of points to your plot • Each geom function takes a mapping argument. The mapping argument is always paired with aes() • ggplot(data= <DATA>)+ GEOM-FUNCTION(mapping=aes(<MAPPINGS>)) • We will learn how to complete and extend this basic template.

> ggplot(data=mpg) +

+ geom_point(mapping=aes(x=displ, y=hwy, colour=class))



Added a third variable called class to a 2-D scatter plot by mapping it to an *a*esthetic.

• Map an aesthetic to a variable —		
• Associate the name of the aesthetic to the name of the variable	••	
 Above example Name=colour and Variable=class. 		
• Scaling: ggplot2() will assign a unique level of the aesthetic col	Lour to a unique level to the variable class.	
 ggplot2 will also add a legend explaining the levels 		
• Other aesthetics include : shape and size		

- The viridis scales provide colour maps that are perceptually uniform in both colour and black-and-white.
- They are also designed to be perceived by viewers with common forms of colour blindness.
- See also https://bids.github.io/colormap/.

ggplot()-viridis options

> ggplot(data=mpg) +

- + geom_point(mapping=aes(x=displ, y=hwy, colour=class))+
- + scale_colour_viridis_d()



Using colour palette from viridis package (colour blind colours).

ggplot()-viridis options



Using colour palette from viridis package (colour blind colours).

- As aesthetics was used to add an additional variable to the plot, another way is to add facets (useful for categorical variables).
- facet_wrap splits plot by a single variable into subplots that each display one subset of the data.

ggplot()-facets



ggplot()-facets and Viridis



- + geom_point(mapping=aes(x=displ, y=hwy, colour=class))+
- + scale_colour_viridis_d(option = "plasma")+
- + facet_wrap(~class, nrow=2)



facet_wrap splits plot by a single variable into subplots that each display one subset of the data.

ggplot()-facets

- > ggplot(data=mpg) +
- + geom_point(mapping=aes(x=displ, y=hwy, colour=class))+
- + scale_colour_viridis_d(option = "plasma")+
- + facet_grid(drv~cyl)



facet_grid splits plot by a combination of two variables into subplots that each display one subset of the data.

ggplot()-geoms

- > ggplot(data = mpg) +
- + geom_point(mapping = aes(x = displ, y = hwy))



a geom is the	
heoretrical object	
lhat a plot uses	
to represent	
data.	

- > ggplot(data = mpg) +
- + geom_smooth(mapping = aes(x = displ, y = hwy))





ggplot()-geoms

ggplot()-geoms









```
> ggplot(data = mpg) +
```

```
+ geom_bar(mapping=aes(x=class,fill=trans))+
```

```
+ scale_fill_viridis_d()
```





```
> ggplot(data = mpg) +
```

+ geom_bar(mapping=aes(x=class,fill=class))+

```
+ scale_fill_viridis_d() + coord_flip()
```



> ggplot(data = mpg) +

- + geom_bar(mapping=aes(x=class,fill=class))+
- + scale_fill_viridis_d() + coord_polar()



class

clas	SS	
	2seater	
	compact	
	midsize	
	minivan	
	pickup	
	subcompact	
	suv	

. Scatter plot	vindis cobuo	
· Barcharts	- facets	
- geons	· coordinate	
:] many		
. more		

```
ggplot(data = <DATA>) +
        <GEOM_FUNCTION>(
            mapping = aes(<MAPPINGS>),
            stat = <STAT>,
            position = <POSITION>
        ) +
        <COORDINATE_FUNCTION> +
        <FACET_FUNCTION>
```