# Data Visualization with R
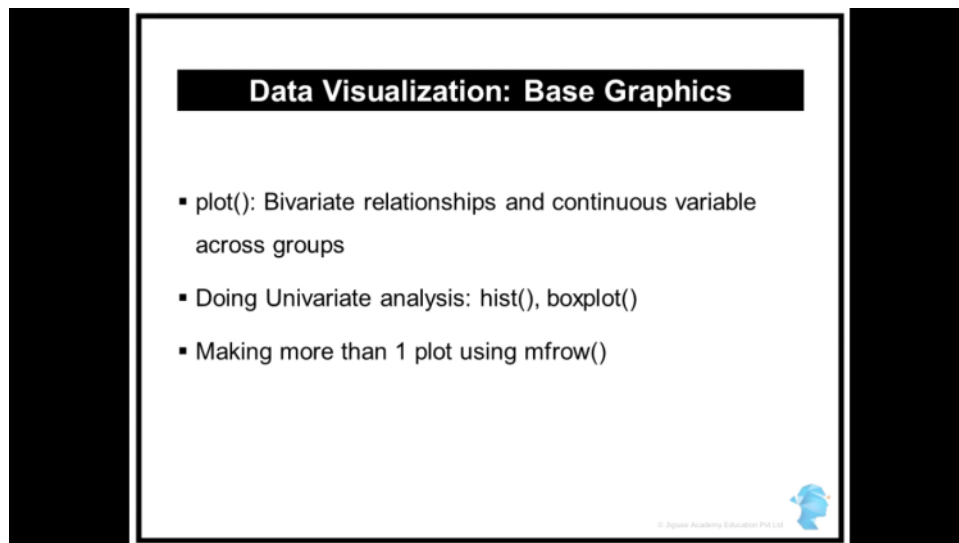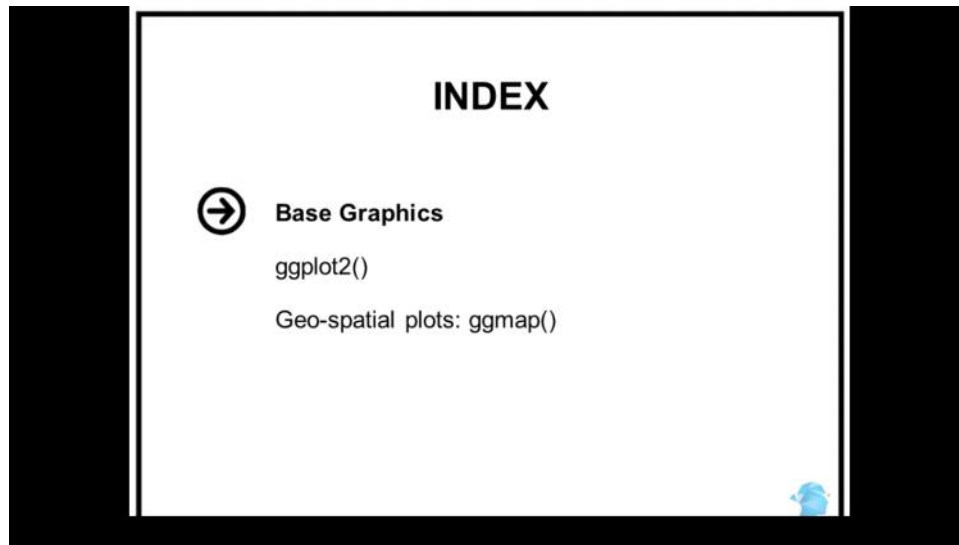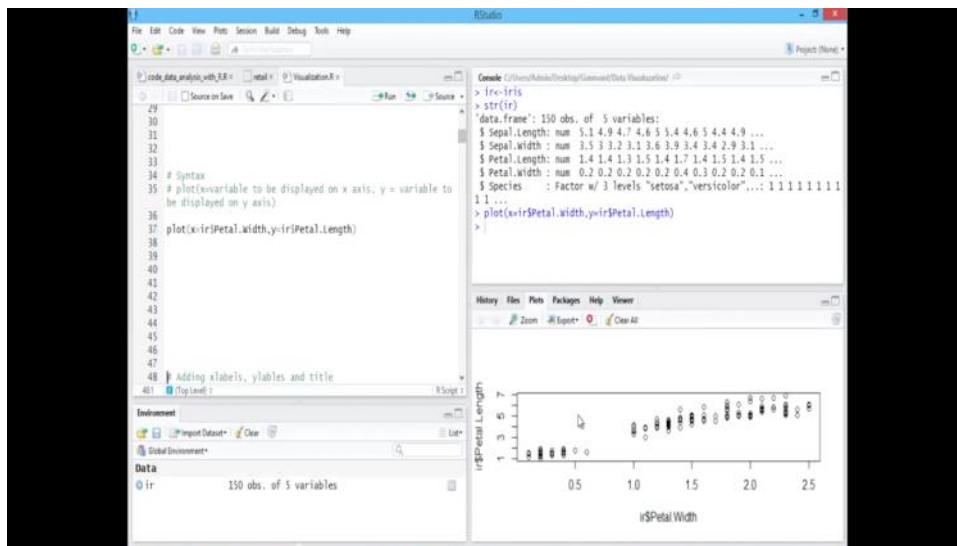
Friday, September 09, 2016      3:21 PM





Lets look at the basic plot function under ggplot2

Plot(x=<variable that we want in the X axis>, y=<variable that we want in the Y axis>)

This is a very basic plot. Now if we want to add title to the plot along with Labels for X and Y axis

Plot(x=iris$petal_length, y=iris$petal_width, main=c("Petal width Vs Length"), xlab=c("Length"),ylab=c("Width"))

To add color to this plot use this:
Plot(x=iris$petal_length, y=iris$petal_width, main=c("Petal width Vs Length"), xlab=c("Length"),ylab=c("Width"), col="red")

If we want to change the default markers from circles to another form, we need to add the PCH attribute to the plot function:

Plot(x=iris$petal_length, y=iris$petal_width, main=c("Petal width Vs Length"), xlab=c("Length"),ylab=c("Width"), col="red", pch=2)
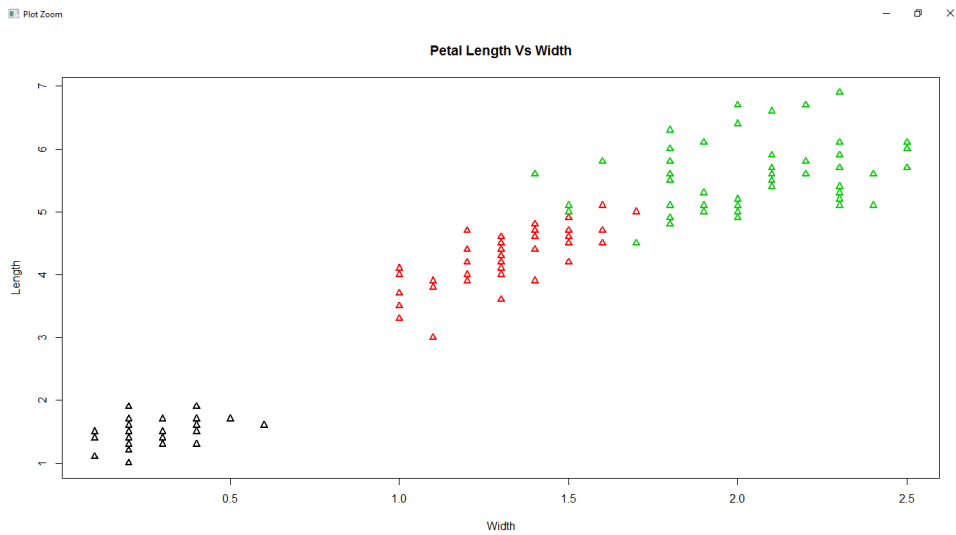
Another optional attribute is the line width parameter:
Plot(x=iris$petal_length, y=iris$petal_width, main=c("Petal width Vs Length"), xlab=c("Length"),ylab=c("Width"), col="red", pch=2, lwd=2)

Conditional Bivariant Plots:

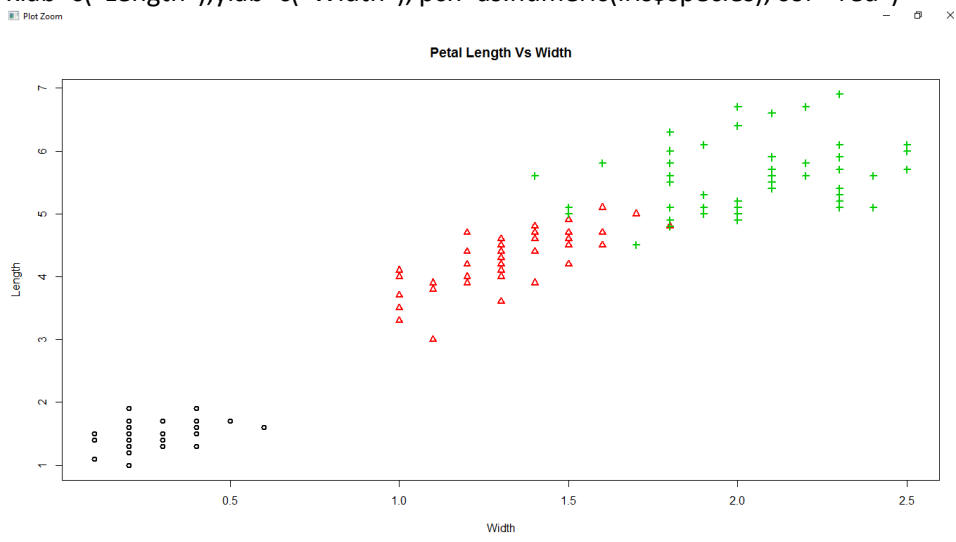We can add conditions to the "col" attribute of the plot function.
This helps in grouping the data in the plot using color coding features that are generated by R.

Plot(x=iris$petal_length, y=iris$petal_width, main=c("Petal width Vs Length"), xlab=c("Length"),ylab=c("Width"), col=iris$species, pch=2)
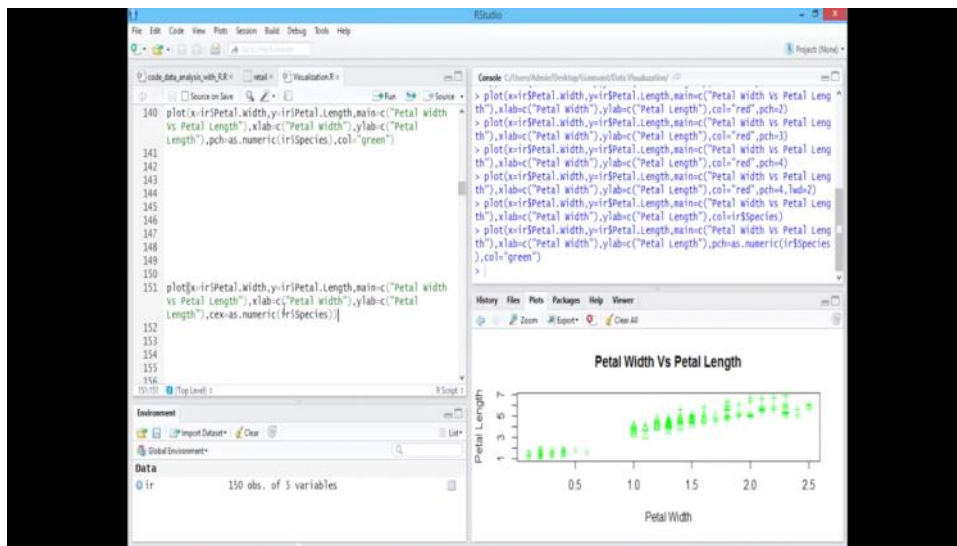
**Petal Length Vs Width**



To do conditional plots based on different chracters , rather than color, we can use the PCH attribute to load the condition: Make sure to use the as.numeric attribute
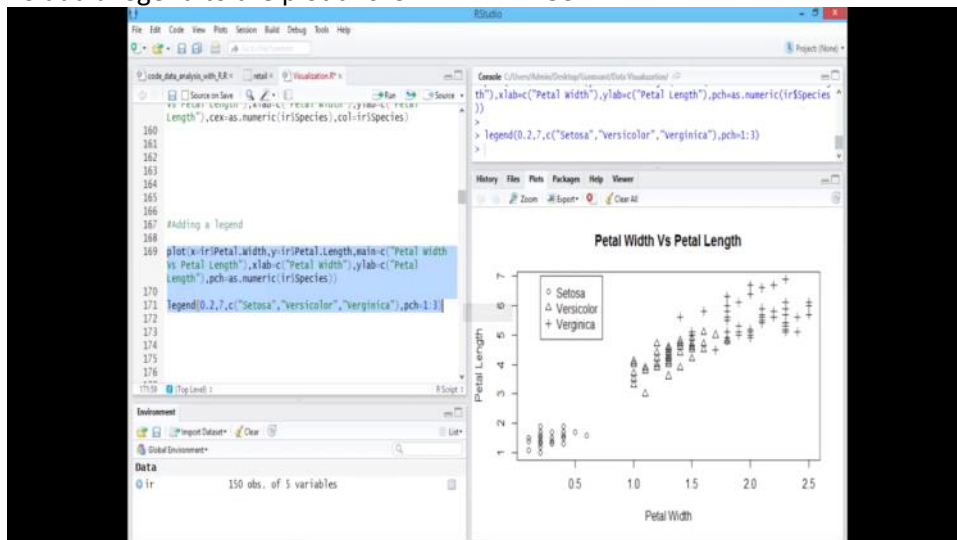
Plot(x=iris$petal_length, y=iris$petal_width, main=c("Petal width Vs Length"), xlab=c("Length"),ylab=c("Width"), pch=as.numeric(iris$species), col="red")
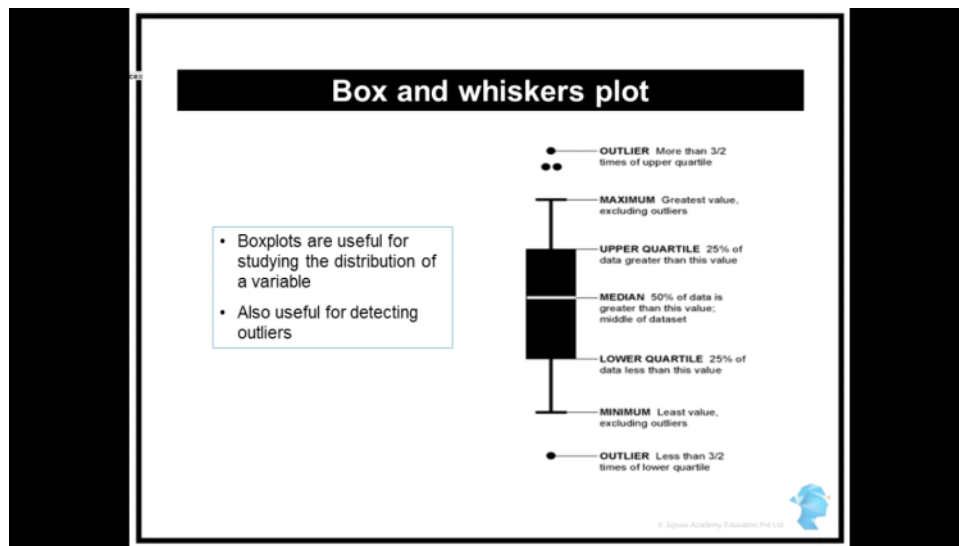
**Petal Length Vs Width**



Another way to differentiate the species is to use the CEX command. This will differentiate based on the size :

To add a legend to the plot this IS WHAT WE USE:





Studying Univariates()

Inorder to run a boxplot we will use the boxplot function:
Boxplot(variable)

Here is how we can use the PLOT function to plot boxplots:
plot(x=iris_test$Species, y=iris_test$Sepal.Length, main = c("Sepal Length across Species"), xlab = "species", ylab = "Length", col="red")

On the X attribute add the factor based on which we want to perform the box plot and Y axis has the attribute for which
We want to perform the analysis.



# Histograms:

To get a histogram plot we can use the HIST function.

hist(cars$dist, labels = TRUE, breaks = 15, xlim = c(0,150), main=c("Stopping Distance"), xlab="Distance in Ft", col="orange")

Here lables=TRUE will show the counts of frequency in the data points.

If we supply the value: freq=FALSE.
This will plot a density function.

If we want to plot a line on this histogram, use the lines command

Lines(density(iris$sepal.width))

In order to display multiple plots on the same area we will use the above function.

PAR --> plotting area function

Mfrow--> number of rows and columns that we want the plotting area to be broken into.

Par(mfrow=c(1,2))

Then run the two distinct plot ffunctions to get the following results:

If we want to return to the default plotting area of just having a single plot, use:
Dev.off()

## Visualization: ggplot2()

- ggplot2(): What and Why
- ggplot2(): Architecture : Understanding Grammar of Graphics
- ggplot2(): Common plots

## Visualization: ggplot2()

- Base graphics: Good for simple tasks
- Comparatively difficult syntax
- Based on grammar of graphics: Simple syntax, interfaces with ggmap and other packages

## Grammar of Graphics

This is the architecture on which GGPLOT2 is built.

## Visualization: ggplot2()

- "Grammar of graphics"
- A plot composed of : Aesthetic Mapping, Geoms, Statistical Transformations, Coordinate Systems and Scales

| Components | Description |
|---|---|
| Aesthetic Mapping | What component of data appears on X axis, Y axis, how is the color, size, fill and position of elements is related with the data |
| Geoms (Geometrical Objects) | What geometrical objects appear on the plot: points, lines, polygons, area, boxplot, rectangle, tile etc |
| Statistical Transformations | Compute density, counts, (Histogram: Need to bin and count data) |
| Scales and Coordinate System | Discreet scale or Continous. Cartesian or Spherical. |

---

## Visualization: ggplot2()



**Aesthetics:**
**Axis Mappings:**
X=temp, y=dewpoint
**Colour:** Seasons

- Based on "grammar of graphics"
- Components: Aesthetics, Geoms and Statistical Transformations

---

This is based on a sample data for seasons, temperature and dewpoints.

## Visualization: ggplot2()



**Geoms:**
**Points (Scatter plot)**
Bars, Lines, Polygons, Area,
Density, Boxplots....

- Based on "grammar of graphics"
- Components: Aesthetics, Geoms and Statistical Transformations

Learning to write code for GGPLOT2



```
p<-ggplot(dataset, aes(x=temp,y=dewpoint, colour=season))
p+geom_point(stat="identity")
```

## Visualization: ggplot2()

- How to code the grammar?
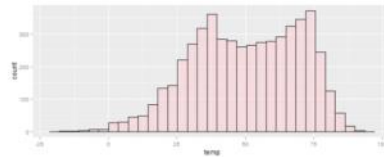


Adding layer with geom:

```
> p+geom_point(stat="identity")
```

- p+geom_type of geometry(stat="statistical transformation on data")

---

## Visualization: ggplot2()



```
> p3<-ggplot(ch,aes(x=temp))
> p3+geom_histogram(colour="black",fill="pink",stat="bin")

> p3<-ggplot(ch,aes(x=temp))
> p3+geom_histogram(colour="black",fill="pink")
```

---

## Visualization: ggplot2()

| Geom | Default Stat | Default Aesthetics |
|---|---|---|
| geom_point | "identity" | colour,fill,shape,size,**x,y** |
| geom_histogram | "bin" | colour,fill,linetype,size,weight,**x** |
| geom_density | "density" | colour,fill,linetype,size,weight,**x,y** |
| geom_polygon | "identity" | colour,fill,linetype,size,**x,y** |
| geom_line | "identity" | colour, linetype, size, **x, y** |
| geom_tile | "identity" | colour, fill, linetype, size, **x, y** |
| geom_boxplot | "boxplot" | colour, fill, lower, middle, size, upper, weight, **x, ymax, ymin** |

*Items in bold are required, others are optional and have default values or are computed by a default stat transform

Lets look at customer demography data:





To further segregate the data based on gender and to add lables to the graph:
We can add the AES function to the geom_point function
i.e. p+geom_point(aes(colour="Gender"))+xlab("salary")+ylab("expenditure")

p+geom_histogram(aes(fill=Species))



p+geom_histogram(aes(fill=Species),position = "dodge")

p+geom_histogram(aes(fill=Species))+facet_grid(Species~.)
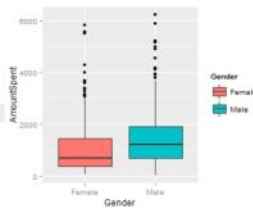


For further polishing the graph:

BOX PLOTS:

Visualization: ggplot2()

- Understanding bivariate counts : 2 d bivariate plots, 2d heatmaps

Geoms:
geom_bin2d()

Aesthetics:
Axis Mappings:
X= X variable
Y= Y variable
Colour, fill......
Statistical Transformation
2 d density

---

# RECAP

- Grammar of graphics
- Code the grammar in R: ggplot2()
- Some basic plots in ggplot2()

---

# Geospatial plots: ggmap()

**Visualization: ggmap()**

- Downloading maps using ggmap()
- Overlaying the map with long-lat data
- Extracting long-lat data from shape files using rgdal() package



**Visualization: ggmap()**

Overview:

Downloading maps ggmap()

↓

Get long-lat data
Text file
Geospatial file: rgdal()

↓

Overlay data on the map:ggplot2()



**Downloading maps: ggmap()**

Here is the command to get the map of a location that we want to plot:

Map<- get_map("North Carolina", maptype="hybrid")

**Visualization: Downloading maps using ggmap()**

Downloading maps using ggmap()

```
> map<-get_map("bangalore",maptype="hybrid")
```

location name

"hybrid",
"terrain",
"satellite",
"roadmap"



**Visualization: Overlaying data on maps**

Overlaying data on maps using ggplot2()

Notice aesthetic mapping

```
> ggmap(map)+geom_point(data=sh,aes(x=long,y=lat),colour="red")
```

ggplot2()

Example:

Ggmap(map)+geom_point(data=sh, aes(x=long, y=lat), colour="red")

SpatialPointsDataframe:

Shape2<- readOGR(dsn="foldername", "filename")
Class(shape2)

This will be of class SpatialPointsDataFrame

To get to know the components of this dataframe we will use the @ symbol and not $
i.e. shapes2@



If the data is in northings and eastings, we will need to convert the coordinates to lat-long format
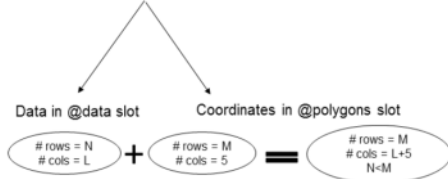
i.e sptransform(shape2, CRS("+init=epsg:4326"))
Once this transformation is done, we will combine the two dataframes

Datac<-data.frame(shape2@data, shape2@coords)

These are polygon dataframes



To extract the lat-long data, we will use the fortify function:



Merge the two datasets based on ID