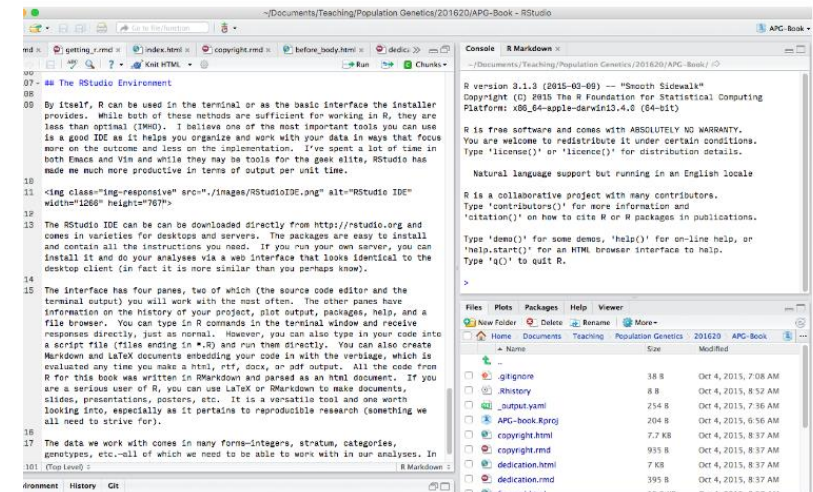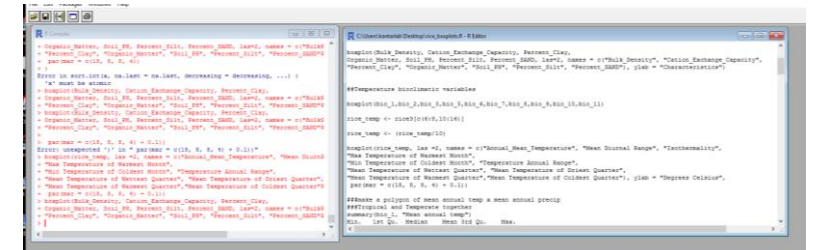# The R ecosystem

# Download 'R'

- R is open source software developed from proprietary software developed by AT&T

- The software works on all platforms and can be downloaded at:
  - https://cran.r-project.org/

- A useful way to use this software is RSTUDIO
  - https://www.rstudio.com/home/

# Getting Started

- There is only a limited amount of functionality in R
  - It will only do what you tell it to do

- People write 'scripts' (software) known as packages that you load to add specific functionality

- It is good to use specific packages to answer your specific research questions, and if there isn't a function you can write one yourself

- You tell R to use a package by using the command
  - library()

# Using R to do basic Math

- Basic operations can be done using
  - Addition: +
  - Subtraction: -
  - Multiplication: *
  - Division: /
  - Exponentiation: ^
  - Modulo: %%
  - sqrt()
  - = assigns value, is not an operation

# Assign a value to a name

rainbow<-5

rainbow<-c(1,2,3,4,5,6)

rainbow/2

mean(rainbow)

# Create a matrix
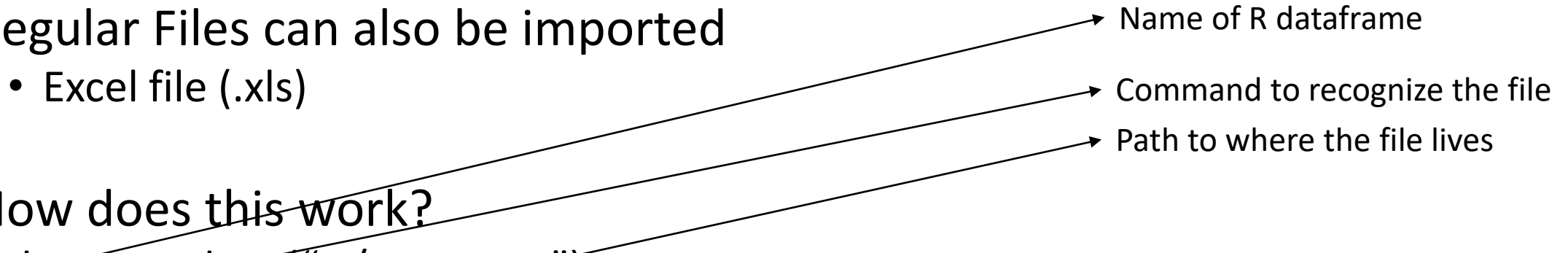
```
mat <- matrix(c(6,5,4,3,2,1),ncol=3,byrow=TRUE)


diag(mat)


mat%*% t(mat)
```

# Inputting Data

- Most common formats are 'flat' files
  - Comma separated files (.csv)
  - Tab delimited files (.txt)
- Regular Files can also be imported
  - Excel file (.xls)

Name of R dataframe

Command to recognize the file

Path to where the file lives

- How does this work?
  data<-read.csv("c:/mycsv.csv")
  data<-read.table("c:/mytext.txt", header=T)
  library(xlsx)
  data <- read.xlsx("c:/myexcel.xlsx", 1)

# Inputting Data

- What if I don't want t write the whole file path

- Print working directory

getwd()

- Set working directory

setwd("dir where your files are")

Print the current directory

Set directory where your files are and where you can call your files, now you just type the file name rather than the entire file path

# Exploring Data: What does the data look like?

#Check if data is correct

names(data) → Gives the headers of all the columns in the data

#look at the data

head(data) → Looks at the first six lines of the dataframe

# Exploring Data: Importance of type of variables

#Check if data is correct

str(data) → Look at what the internal structure of an R object

#look at the data

summary(data) → Generate summary statistics on every variable in your dataset

# Manipulating your data set

- Sometimes you will need to manipulate your dataset

- Remove a column

Subdata<-data[,-1]

- Add a column

data$newfactor<-as.factor(data2$factor)

When exploring a dataframe R recognizes data as [row, column], this command says remove the first column and leave the number of rows the same

Add a new factor from a second dataset

# Manipulating your data set

• Add rownames

rownames(data)<-data[,1]

Take the first column and make these data rownames for each row

• Transpose your dataset

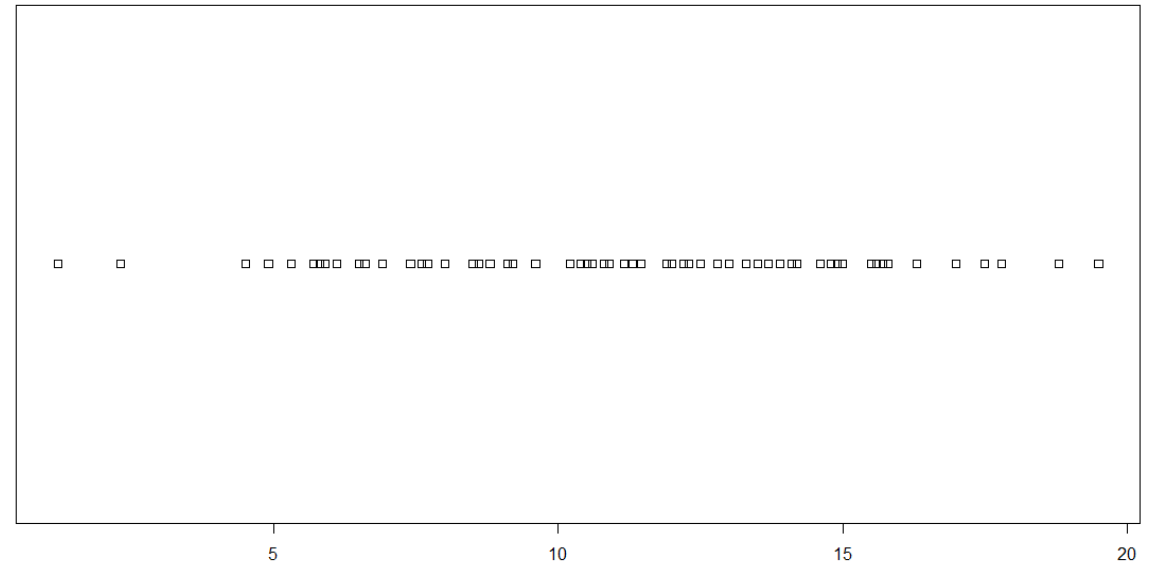Data_transpose<-t(data)

Transpose your matrix

# Visualizing Data: Common Graphics

- Strip Charts

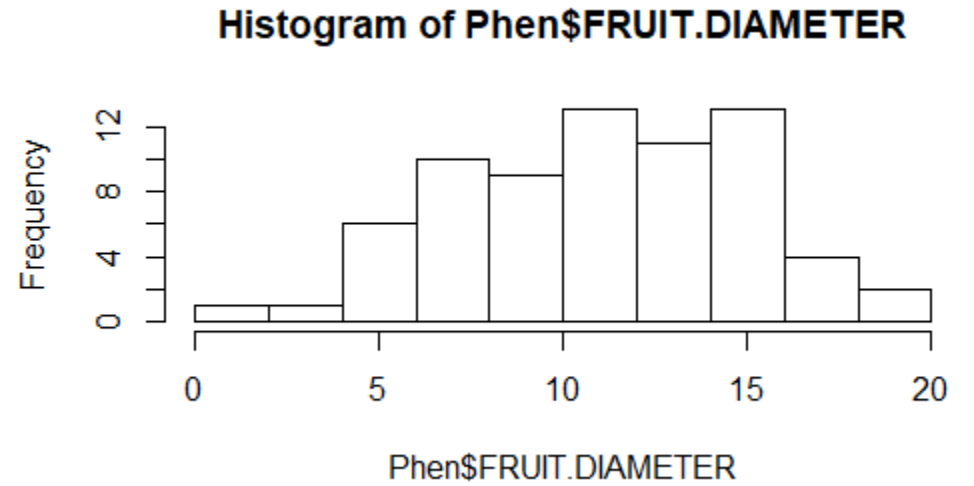- Histograms

- Boxplots

- Scatter Plots

- Normal QQ Plots
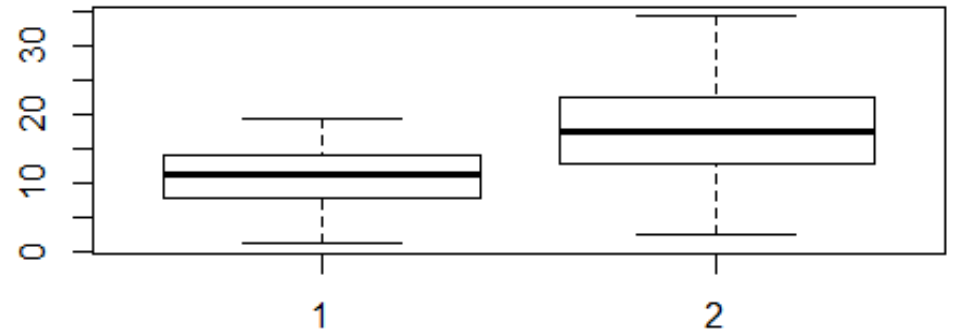
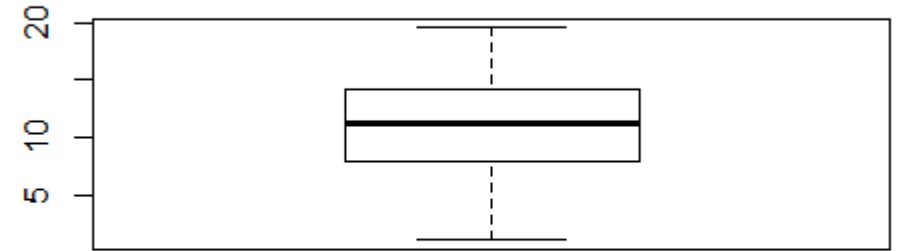# Stripchar

- stripchart(data$factor1)

# Histograms

hist(data$factor1)
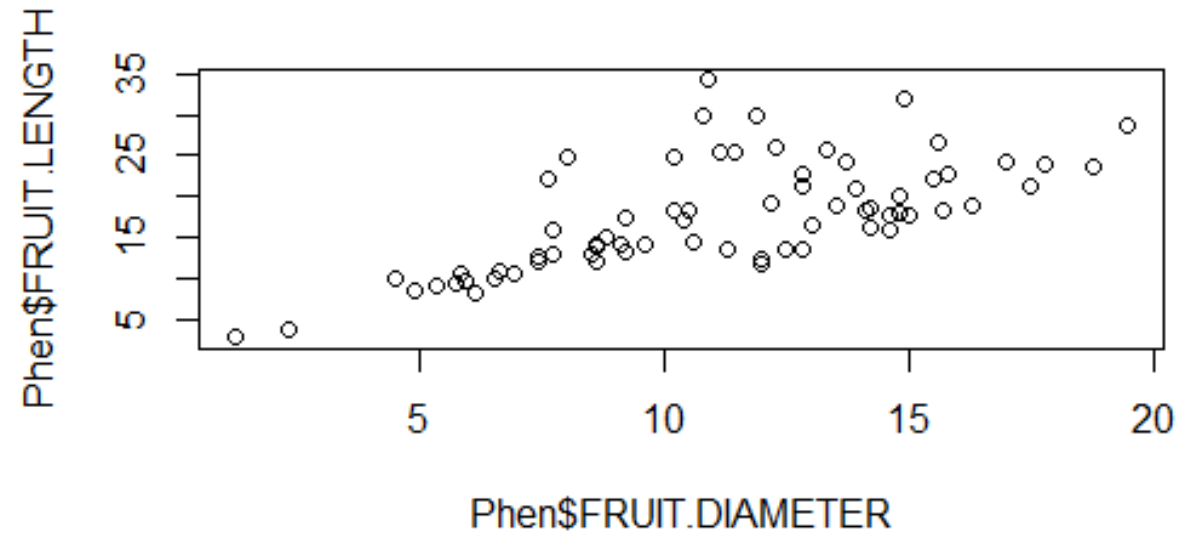


**Histogram of Phen$FRUIT.DIAMETER**

# Box Plots

Boxplot(data$factor1)
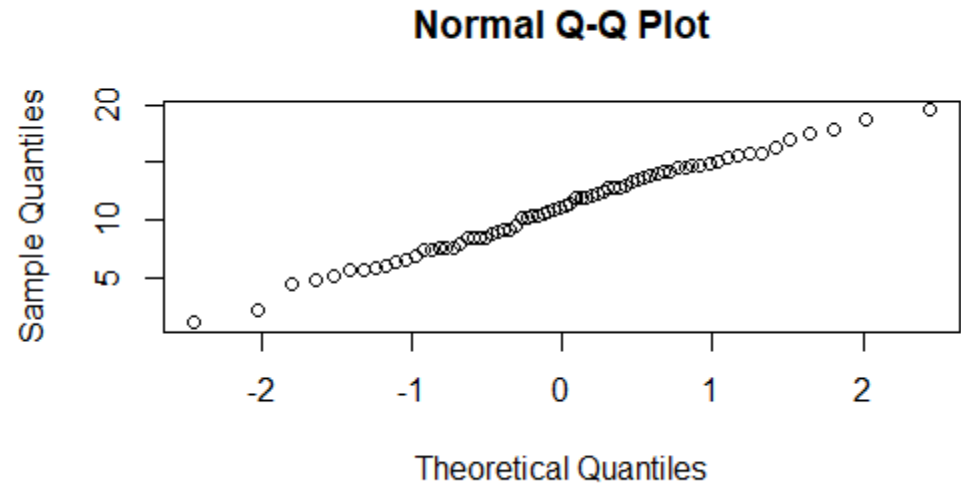


boxplot(data$factor1, data$factor2)

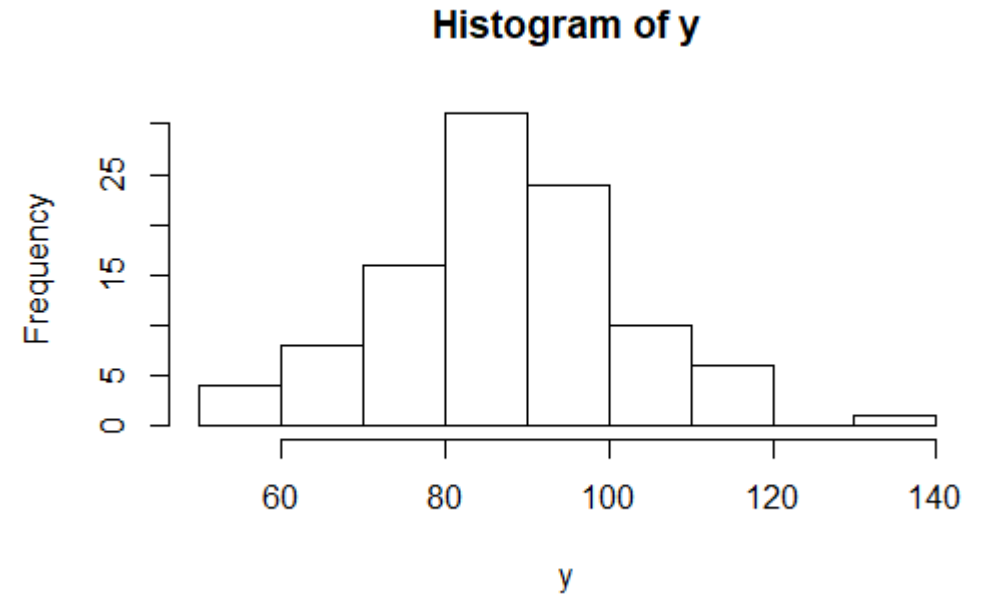# Scatter Plots

plot(data$factor1, data$factor2)

# QQplot

qqnorm(data$factor1,

    main="Normal Q-Q Plot",

    xlab="Theoretical Quantiles",

    ylab="Sample Quantiles")

# Lets Create Distributions

y <- rnorm(100,mean=0,sd=1)
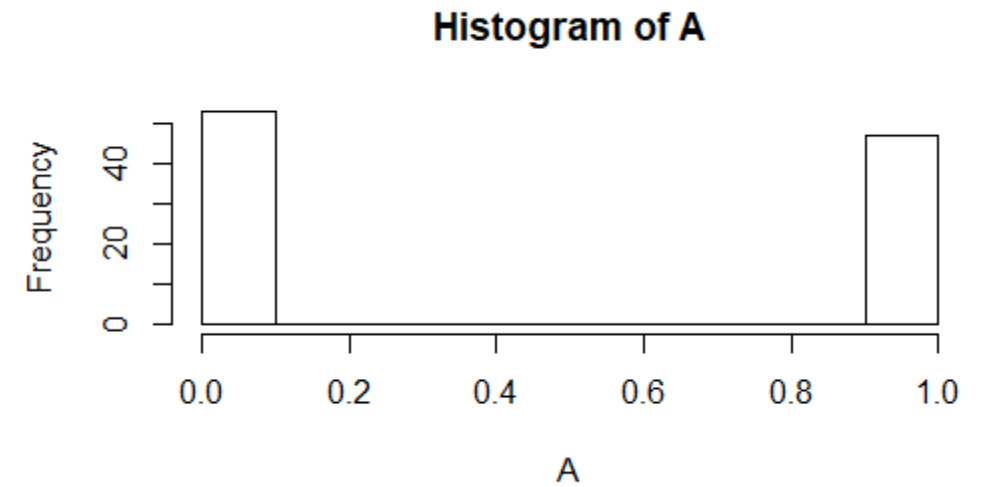hist(y)



Histogram of y

# Bernoulli Distribution

A<-rbinom(100,1,.5)

hist(A)



**Histogram of A**

# The Poisson Distribution

B<-rpois(100, 5)

hist(B)
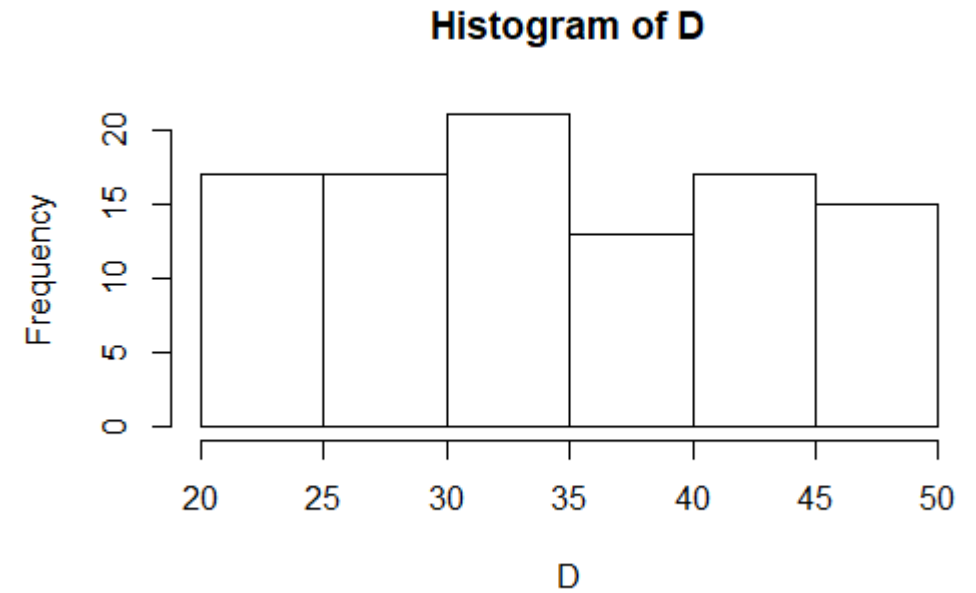


Histogram of B

# Uniform distribution

D<-runif(100,20,50)
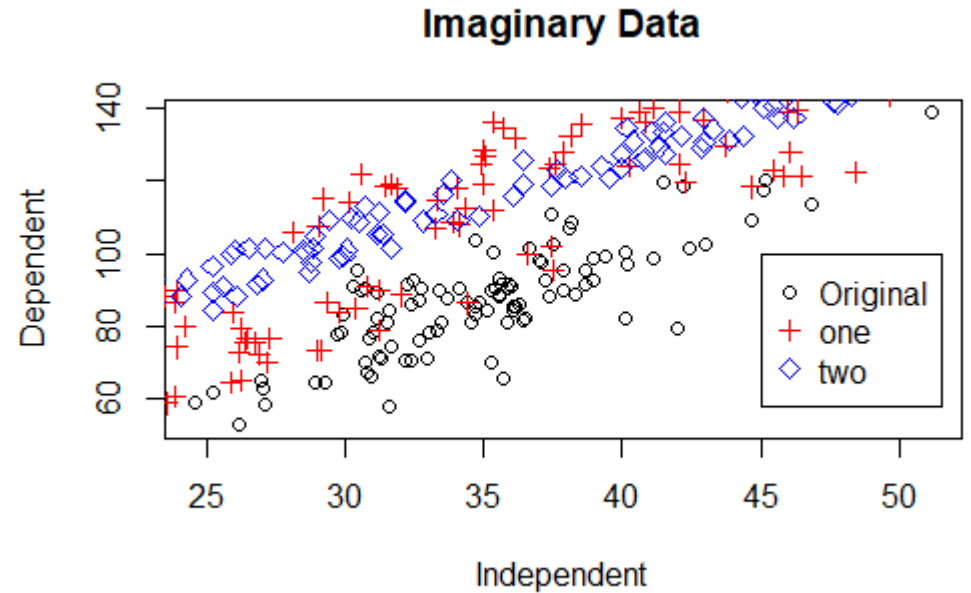
hist(D)

**Histogram of D**

# Let plot some random data

```
x <- rnorm(100,sd=5,mean=35)

y <- 2.5*x - 1.0 + rnorm(100,sd=9,mean=0)

cor(x,y)

plot(x,y,xlab="Independent",ylab="Dependent",main="Imaginary Data")

x1 <- runif(100,20,50) #random uniform distribution, number of samples,min,max

y1 <- 2.5*x1 - 1.0 + runif(100,0,50)

points(x1,y1,col=2)


x2 <- runif(100,20,50)

y2 <- 2.5*x2 - 1.0 + runif(100,22,37)

points(x2,y2,col=3,pch=2)


plot(x,y,xlab="Independent",ylab="Dependent",main="Imaginary Data")

points(x1,y1,col=2,pch=3)

points(x2,y2,col=4,pch=5)

legend(45,100,c("Original","one","two"),col=c(1,2,4),pch=c(1,3,5))
```

# Add error bars

```
x <- rnorm(10,sd=5,mean=20)

y <- rnorm(10,sd=9,mean=10)

plot(x,y,xlab="Independent",ylab="Dependent"
, main="Error bar chart")

xHigh <- x

yHigh <- y + abs(rnorm(10,sd=3.5))

xLow <- x

yLow <- y - abs(rnorm(10,sd=3.1))

arrows(xHigh,yHigh,xLow,yLow,col=2,angle=90,
length=0.1,code=3)
```

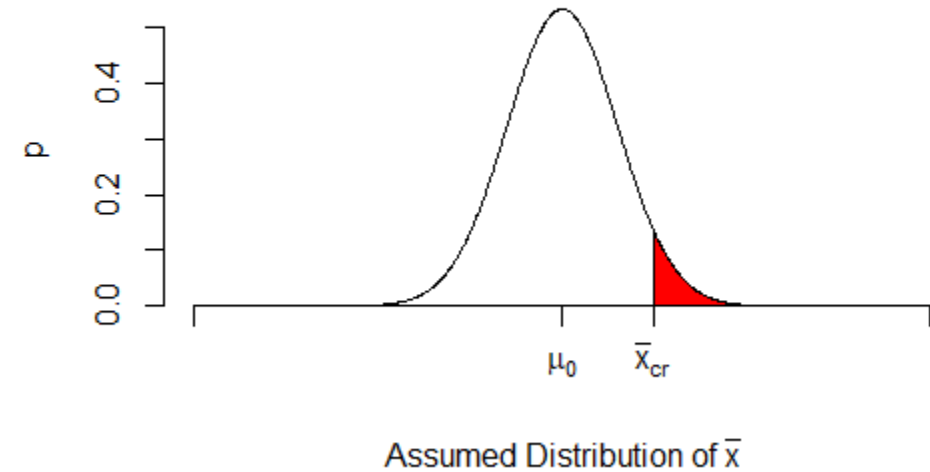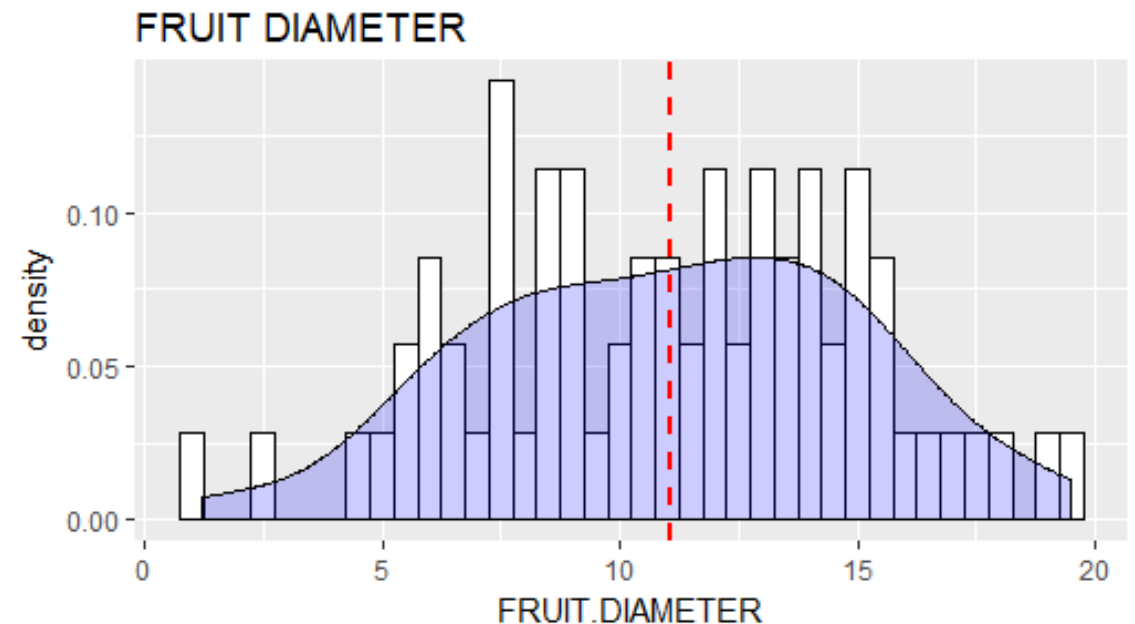# More complicated graphic

```
stdDev <- 0.75;

x <- seq(-5,5,by=0.01)

y <- dnorm(x,sd=stdDev)

right <- qnorm(0.95,sd=stdDev)

plot(x,y,type="l",xaxt="n",ylab="p", # l indicates lines, xaxt removes the x-axis

    xlab=expression(paste('Assumed Distribution of ',bar(x))),

    axes=FALSE,ylim=c(0,max(y)*1.05),xlim=c(min(x),max(x)),

    frame.plot=FALSE)

axis(1,at=c(-5,right,0,5),

    pos = c(0,0),

    labels=c(expression(' '),expression(bar(x)[cr]),expression(mu[0]),expression(' ')))

axis(2)

xReject <- seq(right,5,by=0.01)

yReject <- dnorm(xReject,sd=stdDev)

polygon(c(xReject,xReject[length(xReject)],xReject[1]),

    c(yReject,0, 0), col='red')
```



Assumed Distribution of $\bar{x}$

# More fun plots with ggplot2

library(ggplot2)

ggplot(Phen, aes(x=FRUIT.DIAMETER)) +

  geom_histogram(aes(y=..density..),      # Histogram with density instead of count on y-axis

        binwidth=.5,

        colour="black", fill="white") +

  geom_density(alpha=.2, fill="Blue") + # Overlay with transparent density plot

geom_vline(aes(xintercept=mean(FRUIT.DIAMETER, na.rm=T)),   # Ignore NA values for mean

        color="red", linetype="dashed", size=1)
+ggtitle("FRUIT DIAMETER")

```
ggplot(Food_income_data, aes(x=Percent_food,
y=Percent_health,  colour=as.factor(Year),
group=as.factor(Year), label=Country.Name)) +

geom_smooth(method='lm', formula = y ~ x) +

geom_point() +

geom_text(data=subset(Food_income_data, Year <

2009), size=2, color='black')+

 theme_bw()+

scale_color_brewer(palette = "Set2")+

xlab("Percent of income Spent on Food") +

ylab("Percent of income Spent on Healthcare")+

theme(axis.title.x = element_text(face="bold",
size=12),

axis.title.y = element_text(face="bold", size=12),

legend.title = element_text(size=14, face="bold"))
```

# Another fun graphic

```
##Plot PCA with convex Hulls around the groups

library(ggplot2)

library(plyr)

library(dplyr)

#getting the convex hull of each unique point set

df <- data.frame(pap_1_PCA$ind$coord[,1:5])

df$Kmeans3<-as.factor(fit3$cluster)

df$Kmeans4<-as.factor(fit4$cluster)

df$Kmeans5<-as.factor(fit5$cluster)

df$Pclust<-as.factor(res.hcpc$data.clust$clust)

str(df)

find_hull <- function(df) df[chull(df$Dim.1, df$Dim.2), ]

hulls2 <- ddply(df, "Pclust", find_hull)

#Clustering on 5 groups

plot_of_polygons <- ggplot(data = df, aes(x = Dim.1, y = Dim.2, colour=Pclust, fill = Pclust)) +

  geom_text(aes(label=rownames(Phen)),size = 5,hjust=0,vjust=0)+

  geom_polygon(data = hulls2, alpha = 0.5) +

  geom_point( x=-4.21,y=-1.86 ,color='black',size=3,shape=8) +

  geom_point( x=-1.24, y=3.15, color='black',size=3,shape=8) +

  geom_point(x= 3.35,y=-1.13,color='black',size=3, shape=8) +

  labs(x = "PC1", y = "PC2")
```