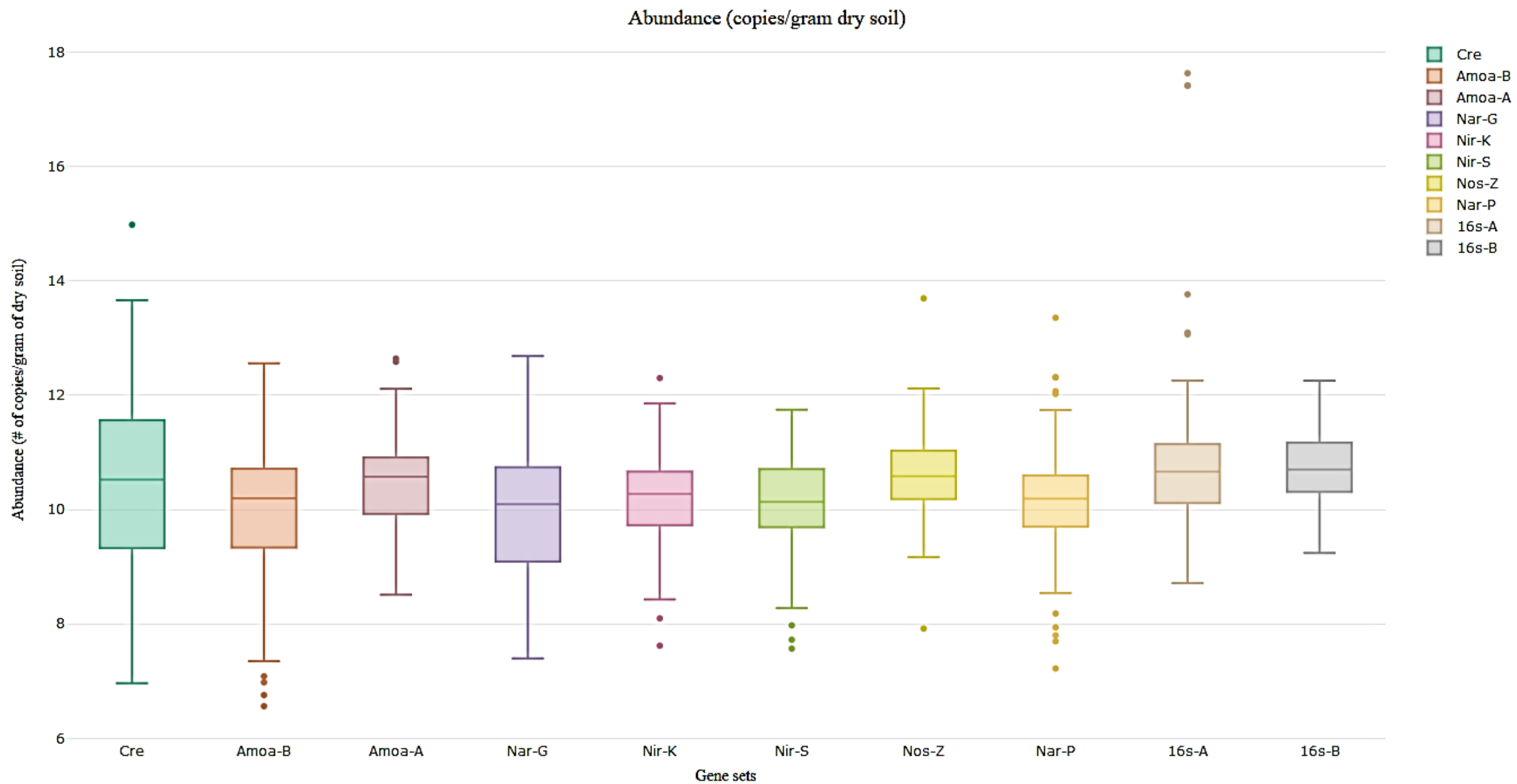# Introduction to R
# and
# Executing General Linear Models

Module 3 & 4

Victor Valdez

# Why use R?

- FREE

  - Powerful and flexible (kind of) statistical and graphical package

- Most frequently used statistical package

  - An increasing number of users with R

  - Open-source software

  - Versatility with coding allows for robust analysis of data

Abundance (copies/gram dry soil)

# What am I looking at?



- **Upper left: the source**
  - **Text editor to save, write and edit code**

- **Upper right: environment tab**
  - **Import dataset**
  - **Lists of objects made**

- **Lower left: the console**
  - **Accepted command lines**

- **Lower right: files tab**
  - **Packages installed**
  - **Images and plots**

# Getting started

- RStudio

  - Makes R easier to use

  - Start by keeping your work organized by:
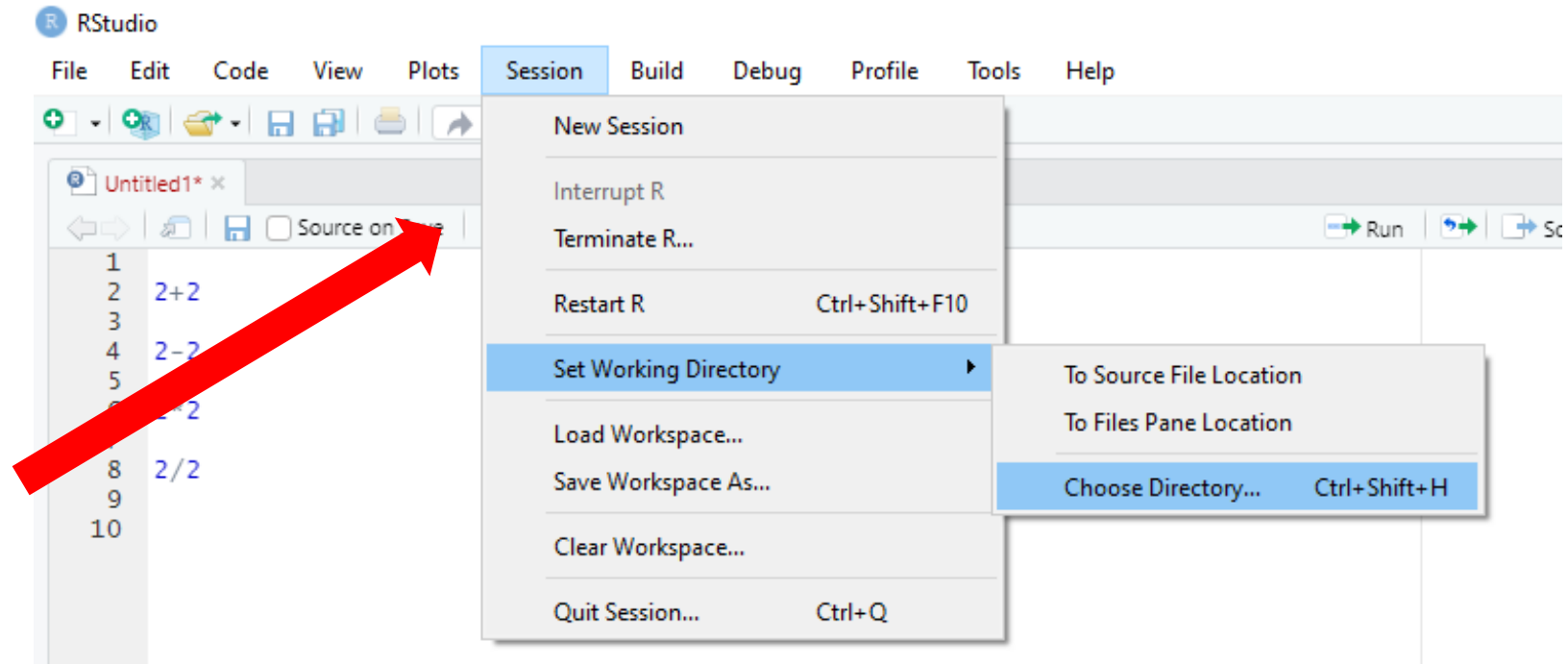
    **Set your *"Working Directory"***

    **OR**

    Ctrl+shift+H

    **OR**

    setwd("C:/Users/Home/Downloads/Workshop_data")

# Getting started

- Write some simple stuff in the text editor

    2 **+** 2

    2 **-** 2

    2 **\*** 2

    2 **/** 2

- **Highlight** code and press "ctrl + enter" or "Run"

- Output is in the console!

# Manual data entry

## Vectors and Matrices

- Create a **vector object** called "a" and "b" of 3 lengths each

    a <-c(0,1,2)

    b <-c(3:5)

- **Type name** and **run code _OR_ click** on object **in global env.**

- Output in console!

- What it means:

    **a** and **b** is the vector object created
    **<-** is the direction items are going
    **c()** is the concatenate function
    **0,1,2** are the items being added _OR_
    **3:5** is an inclusive function

# Manual data entry

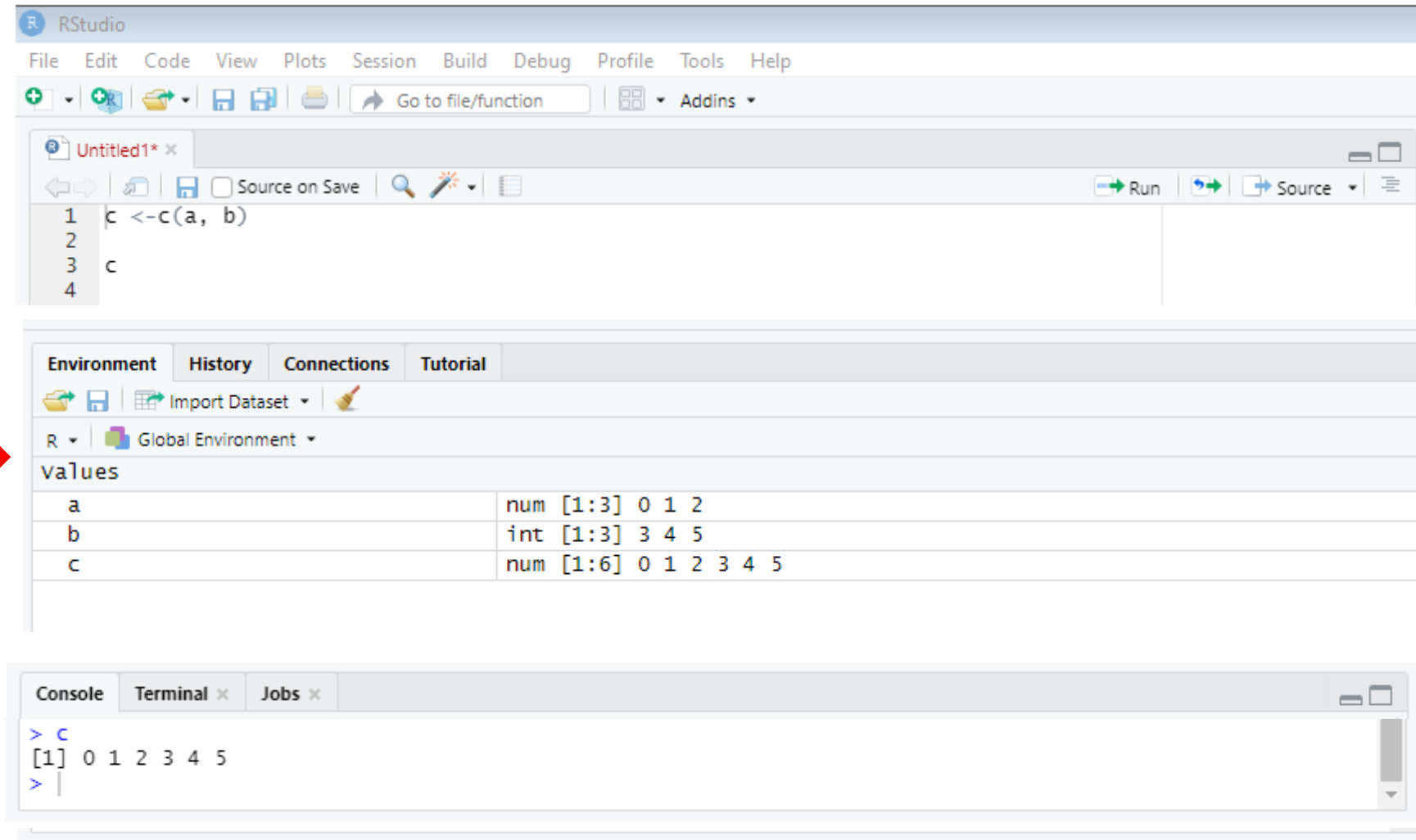## Vectors and Matrices

- Create an object matrix "c" with object vectors "a" and "b"

  c <-c(a, b)

- **Type name** and **run code _OR_ click** on object **in global env.**

- What it means:

  **c** is the new object with **a** and **b**

  **<-** is the direction items are going

  **c()** is the concatenate function

  **a** and **b** are the items being added

# Manual data entry

## Vectors and Matrices
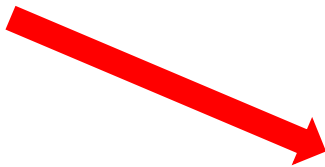
- Create a character vector with letters

    d<-letters [1:10]

- Then bind to numeric vector by column

    cbind(a, d)
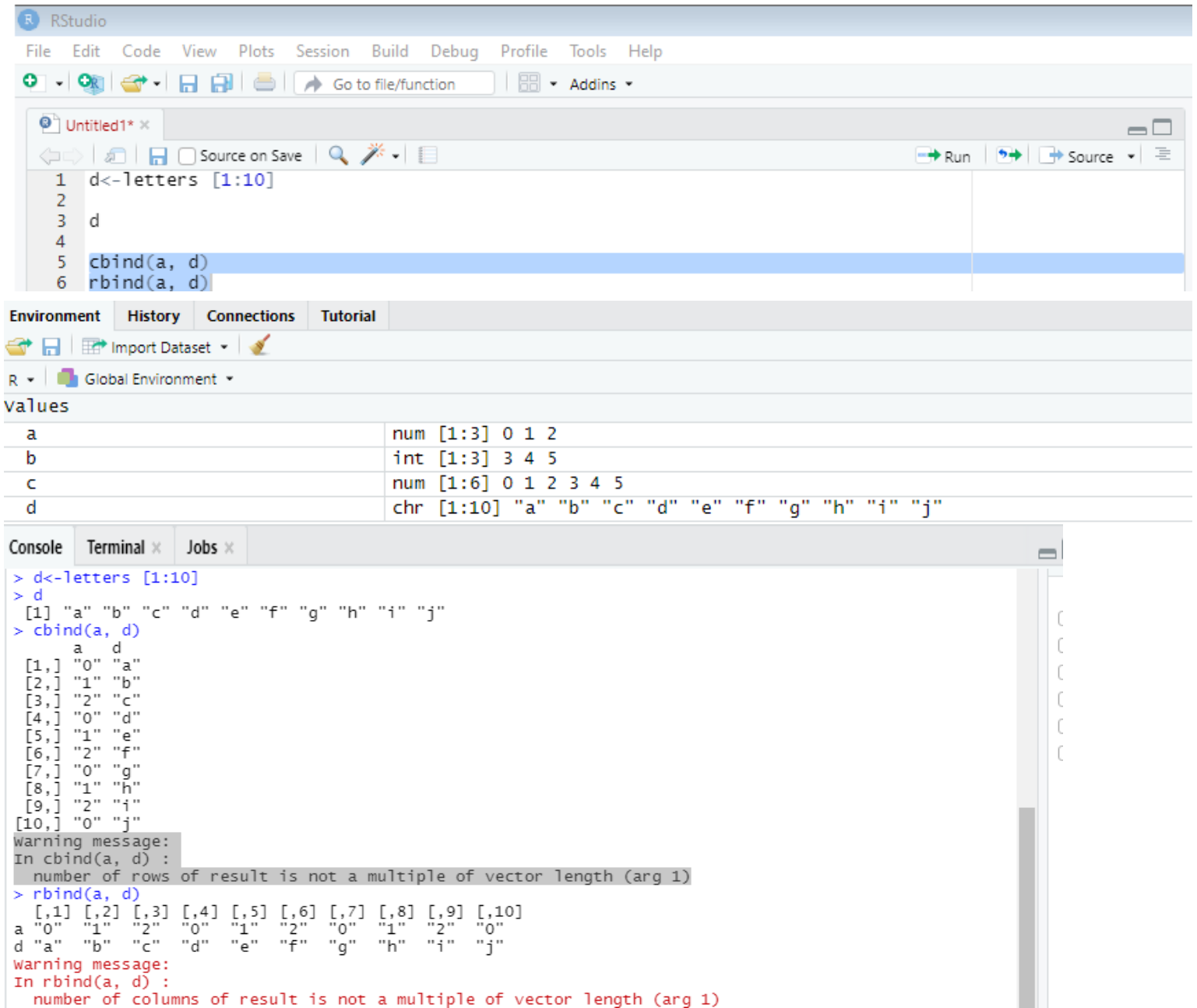
- Now bind by row

    rbind(a, d)

- What it means:

    **Letter []** is the alphabet function
    **[]** specify parts of the object
    **cbind()** bind by column
    **rbind()** bind by row
    **Warning message** R tells us when
    something is wrong. We can ignore this

# Manual data entry

- Create matrix object

    newmat <- matrix(1:9, nrow=3, ncol=3)

    newmat

    newmat[2,3]

- Notice how rows come first then columns

- What it all means:

    **newmat** is a matrix object
    **matrix()** is the matrix object function
    **nrow** is the number of row function
    **ncol** is the number of columns function
    **=** denotes the value to use

# Manual data entry

- To rename any column in your object

  names(newmat)<-c("newname1", "newname2","newname2")



- What it all means:

  **names()** is the renaming function
  **<-** where to put the items
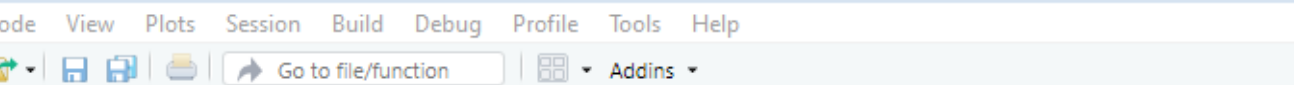  **newmat** is the matrix object
  **c()** is concatenate function
  **newname1** is the new name of the column
  **""** denotes what to rename the columns as

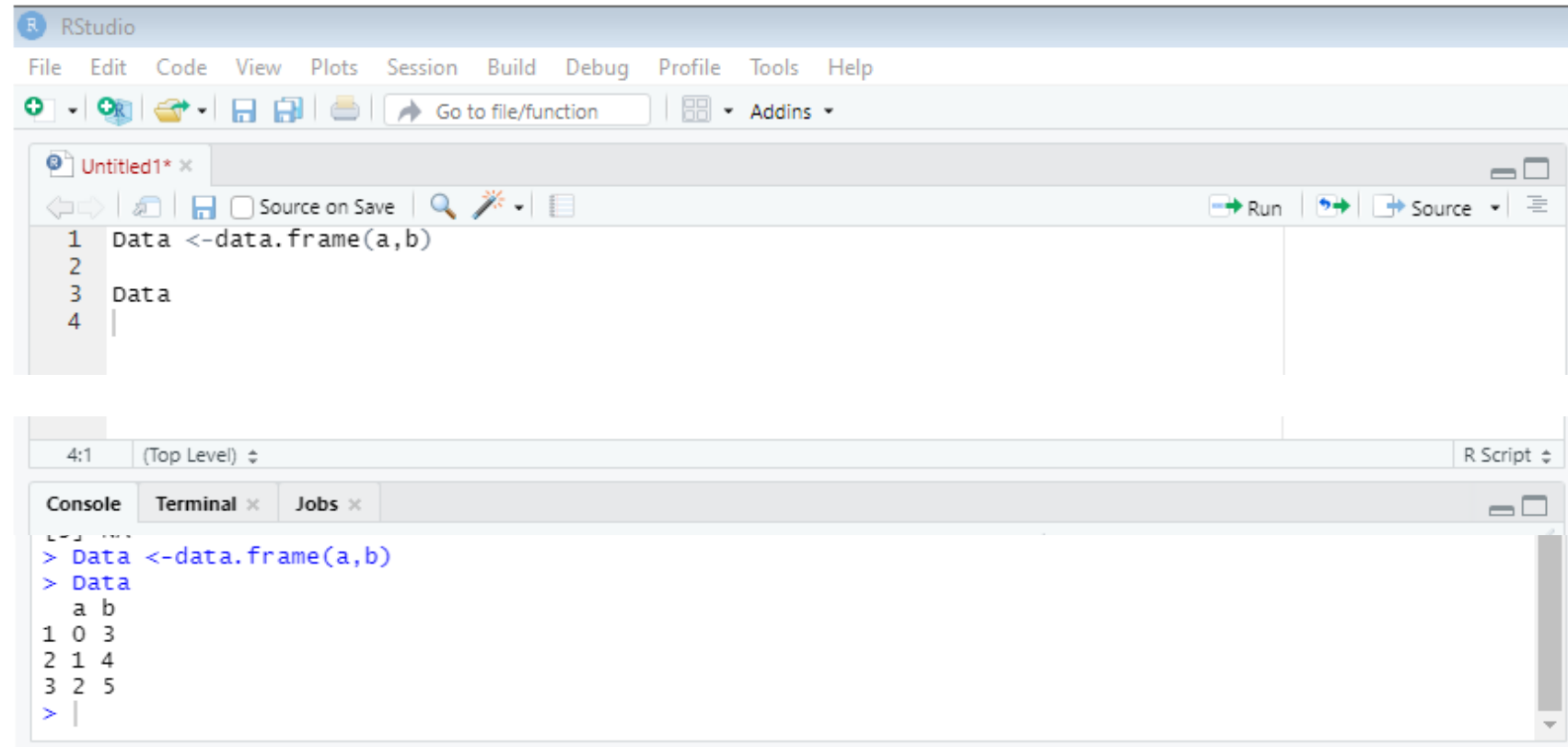  **,** is used to separate the different names

# Manual data entry:

- Make a data frame

  Data <-data.frame(a,b)

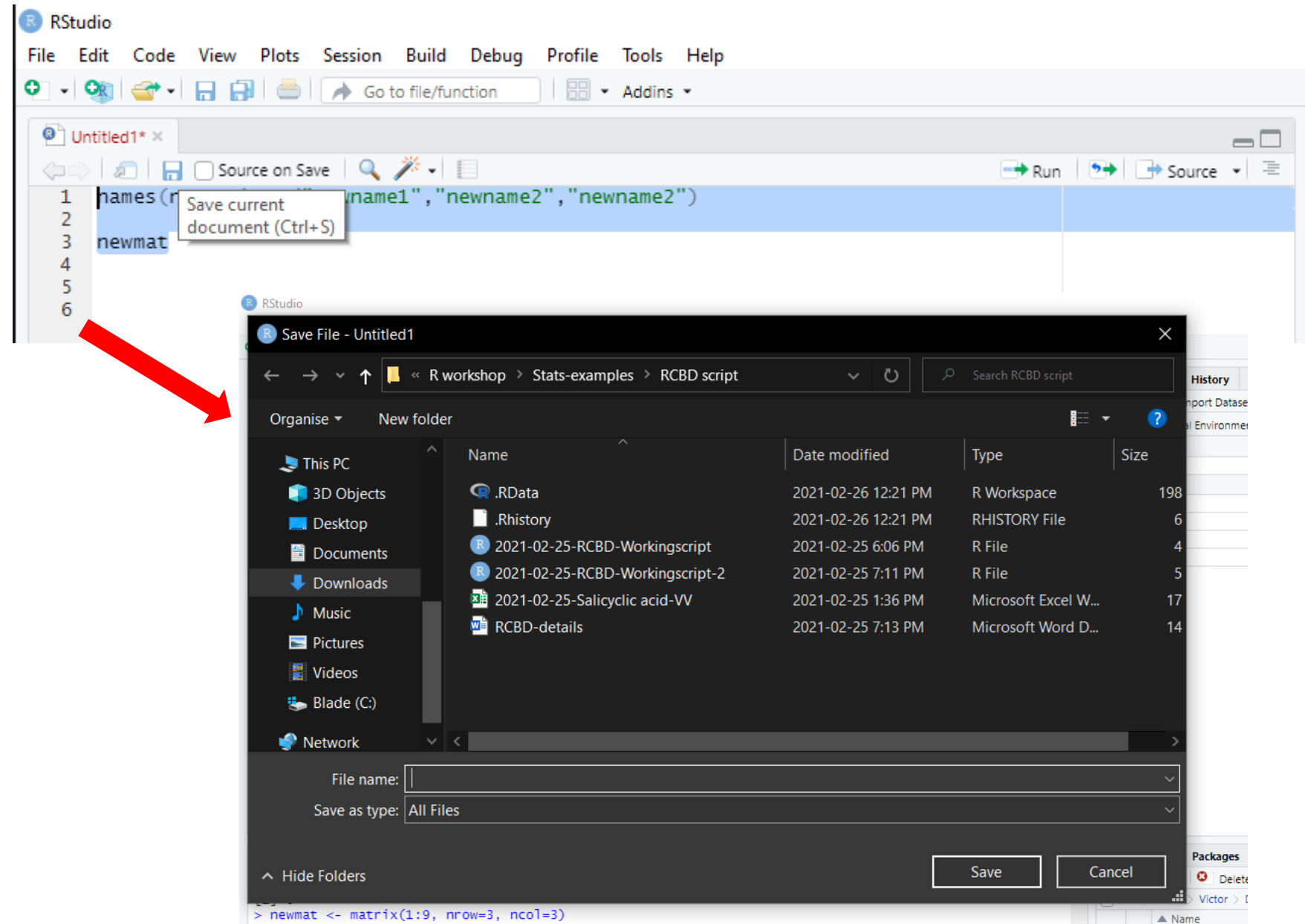- A data frame is similar to an excel sheet; it is a list of vector of equal lengths

- But stores data as a table, can contain multiple data types in multiple columns called fields

- What it means:

  **data.frame ()** is the function that turns objects into a data frame
  **a,b** are your objects

# Manual data entry

- To keep any working code you generated, and to edit it at a later time

- Choose the location of your files

# Descriptive statistics

- Other mathematical functions we can use w/ our data frames

    min (dataset$variable_name)
    max (dataset$ variable_name)
    IQR (dataset$ variable_name)
    mean (dataset$ variable_name)
    median (dataset$ variable_name)
    var (dataset$ variable_name)
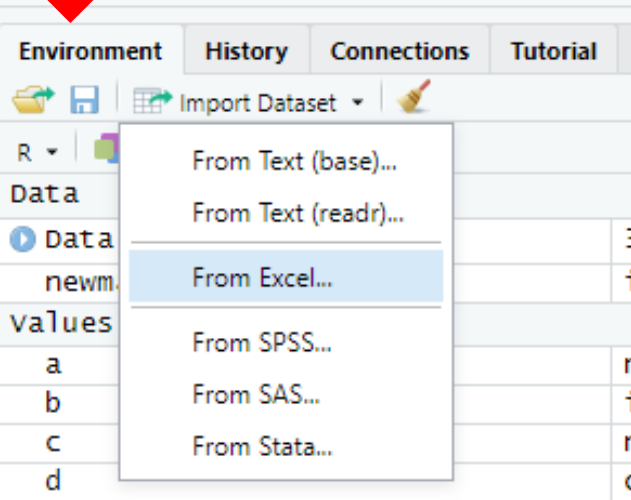    sd (dataset$ variable_name)

- What it means:

    **dataset** is the data frame you are referring to
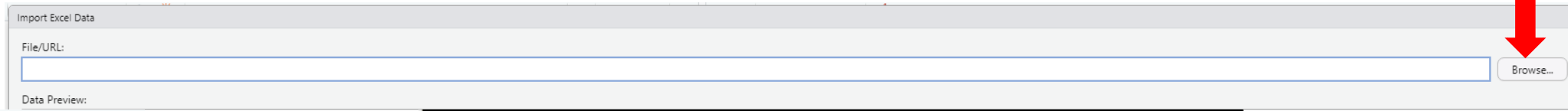    **$** is used to call up the list of columns in the data frame
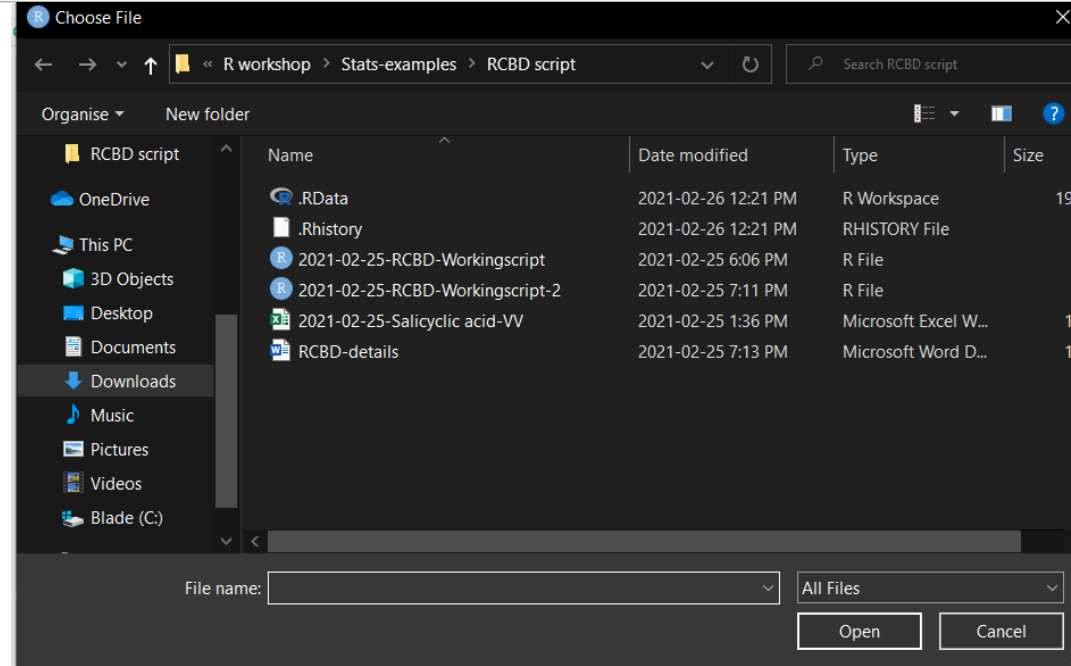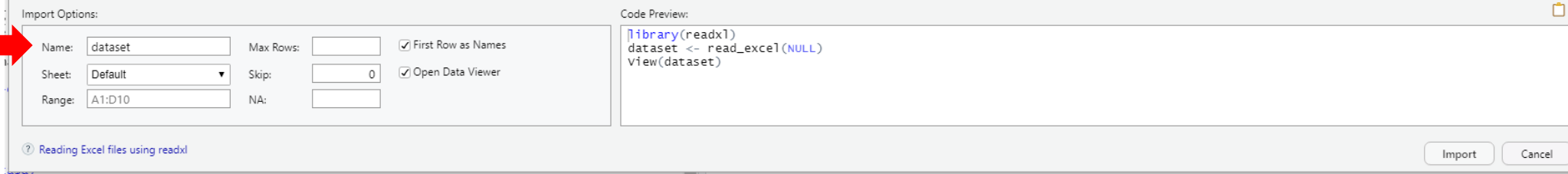    **variable_name** is the column

# Import data into R

- You should now see the following; *__click__* import



- Our excel sheet with columns and data

- We now see the contents of the file name

- The file name may be too long; not ideal for coding

- Let's rename:

sa_data<-X2021_02_25_Salicyclic_acid_VV





- We can now use the autofill feature by using **$ before the dataset**

  e.g., data$....

# Rename Columns and Run Descriptive Statistics

names(sa_data)<-c("b","trt","leaves")

min(sa_data$leaves)
max(sa_data$leaves)
IQR(sa_data$leaves)
mean(sa_data$leaves)
median(sa_data$leaves)
var(sa_data$leaves)
sd(sa_data$leaves)

# Run Descriptive Statistics Function

Change blocks and treatments to factorial:

sa_data$b = factor(sa_data$b)
sa_data$trt = factor(sa_data$trt)

Now use the summary function in FSA package:

sum=Summarize(leaves ~ b, data = sa_data)
sum2=Summarize(leaves ~ trt, data = sa_data)
sum3=Summarize(leaves ~ b+trt, data = sa_data)

Let's look at the tables generated

# Plotting data: extract data

- Data is best plotted, let's make some boxplots

- First extract the mean from the table:

Table = as.table(sum$mean)

rownames(Table) = sum$b

Table2 = as.table(sum2$mean)

rownames(Table2) = sum2$trt

Table3 = as.table(sum3$mean)

# Plotting data

```
63
64  #It's hard to see differences in table form, let's visualize the data:
65  ##Using boxplots()
66
67  B<-boxplot(leaves ~ b, data = sa_data, col = "lightgray")
68  TRT<-boxplot(leaves ~ trt, data = sa_data, col = "lightgray")
69
70  ##there should be no interaction but we look at it anyway
71  ###NOTE:Order matters! try swtiching the order of b and trt
72
73  BTRT<-boxplot(leaves ~ b*trt, data = sa_data, col = "lightgray")
74  TRTB<-boxplot(leaves ~ trt*b, data = sa_data, col = "lightgray")
```

- Plot the boxplot :

B<-boxplot(leaves ~ b, data = sa_data, col = "lightgray")

TRT<-boxplot(leaves ~ trt, data = sa_data, col = "lightgray")

BTRT<-boxplot(leaves ~ b*trt, data = sa_data, col = "lightgray")

TRTB<-boxplot(leaves ~ trt*b, data = sa_data, col = "lightgray")

- Let's view it in the plot output:

# Plotting data

- Plot barplots:

```
barplot(Table, ylab="Mean leaves", xlab="Blocks")
barplot(Table2, ylab="Mean leaves", xlab="Treatments")
barplot(Table3, ylab="Mean leaves", xlab="Blocks and Treatment", co
= Table3)
```

```
75
76    ##Using barplot()
77
78    barplot(Table, ylab="Mean leaves", xlab="Blocks")
79    barplot(Table2, ylab="Mean leaves", xlab="Treatments")
80    barplot(Table3, ylab="Mean leaves", xlab="Blocks and Treatment", col = Table3)
81
```

# Executing assumptions of normality: LMM

- We can execute our command for linear mixed models in two ways in R

   model1 <- lm(leaves ~ b+trt,data = sa_data)

   ***OR***

   model1 <- lm(data$Y~data$X)

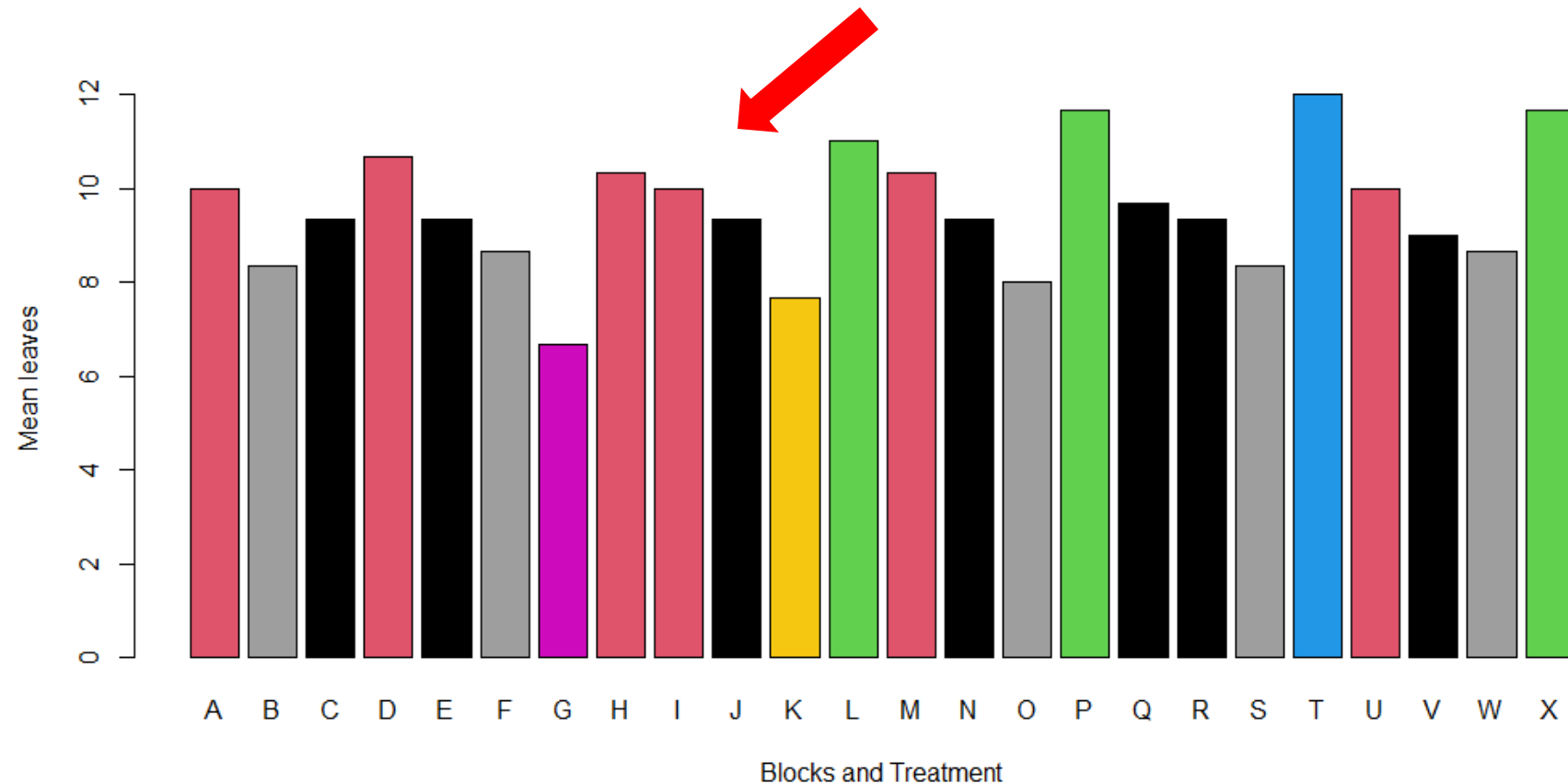- From here we would look at the error structure using the following:

   Res<-resid(model1)

   Fit<-fitted(model1)

   Plot(X,Y, xlab="Fits", ylab = "Residuals", main = "Residual vs. Fits plot for normality")
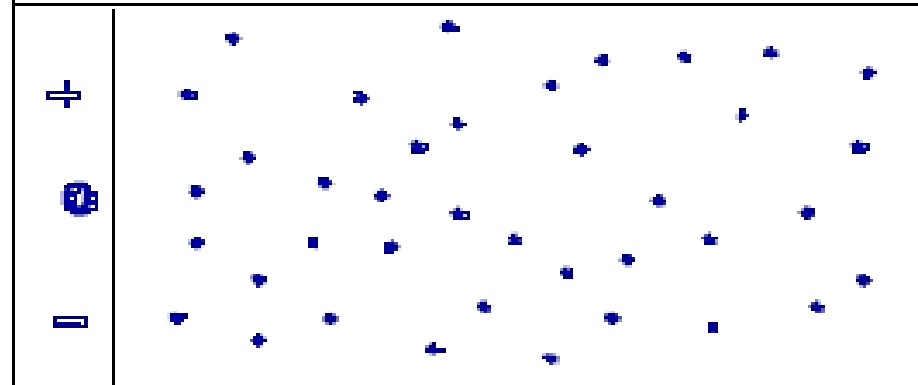
   # Where X = Fit
   # And Y = Res

   abline(h=0, col="red")
   #helps with visualization

# Normal error model assumptions graphs



✓ ACCEPTABLE

✗ NOT ACCEPTABLE

✗ NOT ACCEPTABLE

# Normal error model assumptions: LMM

Homogeneous errors

Use the plot we made

✓ ACCEPTABLE

× NOT ACCEPTABLE

**Residual vs. Fits plot for normality**

# Normal error assumptions: LMM

Independent or autocorrelation errors

## Plot a lag plot

lag.plot(Res, do.lines=FALSE, diag.=FALSE)

✓ACCEPTABLE

× NOT ACCEPTABLE

# Normal error model assumptions: LMM

Normal errors

✓ACCEPTABLE

## Plot a histogram:

hist(Res)



Histogram of Res

× NOT ACCEPTABLE

# Normal error model assumptions: LMM

Normal errors

Make a quantile-quantile normality plot

qqnorm(Res)

qqline(Res)

✓ ACCEPTABLE

× NOT ACCEPTABLE

# LMM ANOVA output

Code and run the ANOVA function

anova(model4)

```
> anova(model1)
Analysis of Variance Table

Response: leaves
          Df  Sum Sq Mean Sq F value     Pr(>F)
b          3  95.111  31.704 18.1484 1.339e-08 ***
trt        5  10.611   2.122  1.2148    0.3126
Residuals 63 110.056   1.747
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
>
```

# LMM ANOVA output

Our blocks can affect the adjusted p-valye let's check the model again with different ordering

model5 <- lm(leaves ~ trt+b ,data = sa_data)

anova(model5)

```
> model2 <- lm(leaves ~ trt+b ,data = sa_data)
> anova(model2)
Analysis of Variance Table

Response: leaves
          Df  Sum Sq Mean Sq  F value    Pr(>F)
trt        5  10.611   2.122   1.2148    0.3126
b          3  95.111  31.704  18.1484 1.339e-08 ***
Residuals 63 110.056   1.747
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
>
```

# LMM ANOVA output

- The block and treatment interaction term can affect the ANOVA
- In this case, the interaction term has a very small weight, let's check

    model6 <- lm(leaves ~ trt*b ,data = sa_data)
    anova(model6)

```
> model3 <- lm(leaves ~ trt*b ,data = sa_data)
> anova(model3)
Analysis of Variance Table

Response: leaves
          Df Sum Sq Mean Sq F value    Pr(>F)
trt        5 10.611   2.122  1.0466    0.4015
b          3 95.111  31.704 15.6347 3.161e-07 ***
trt:b     15 12.722   0.848  0.4183    0.9663
Residuals 48 97.333   2.028
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
> |
```

# LMM ANOVA output

Let's try another approach

library(lme4)

model7 <- lmer(leaves ~ (1|b) + trt, data = sa_data)

anova(model7)

1- pf(1.2148,5,63)

```
> model4 <- lmer(leaves ~ (1|b) + trt, data = sa_data)
> anova(model4)
Analysis of Variance Table
     npar Sum Sq Mean Sq F value
trt     5 10.611  2.1222  1.2148
>
```

# Publication quality boxplots

library(ggplot2)

p1<-ggplot(sa_data, aes(x=trt, y=leaves, fill=trt)) + geom_boxplot()

p2<-ggplot(sa_data, aes(y=leaves, fill=b)) + geom_boxplot() + facet_wrap(~trt,ncol = 6)
+  theme( axis.text.x = element_blank(), axis.ticks.x = element_blank())