

Exploring data - graphical summaries

Introduction to Quantitative Ecology

Fall 2018

Chris Sutherland

csutherland@umass.edu

Group evaluations

1. Which of the following lines of code is the correct way to read a csv file?

- A) `read.csv("mydata.csv")`
- B) `my.data <- read.csv("mydata.csv")`
- C) `r my.data <- read.csv(mydata.csv)`
- D) `read.csv(mydata.csv, row.names = 1)`

2. Which of the following symbols is used to represent the *mean of a variable*?

A) σ

B) \bar{x}

C) s^2

D) x

3. Which plot would be most appropriate for visualizing the *running mean*?
- A) box-whisker plot
 - B) line graph
 - C) histogram
 - D) bar chart

4. Which plot would be most appropriate for visualizing the relationship between *two continuous variables*?
- A) histogram
 - B) line graph
 - C) box-whisker plot
 - D) scatter plot

5. Which of the following plots would I use to graphically represent the *distribution* of a variable?
- A) histogram
 - B) line graph
 - C) box-whisker plot
 - D) scatter plot

Graphical exploration

Why use graphs?

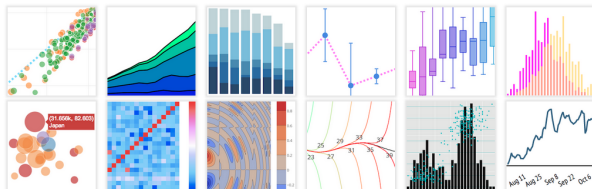


Graphical exploration

Two main reasons to use graphs:

1. Inform how to analyze the data

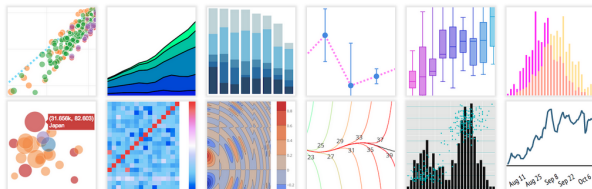
2. Presentation of the data



Graphical exploration

Two main reasons to use graphs:

1. Inform how to analyze the data
 - ▶ visualization
 - ▶ identify patterns
 - ▶ choose appropriate statistical test
2. Presentation of the data



Graphical exploration

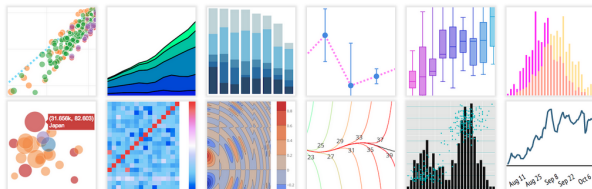
Two main reasons to use graphs:

1. Inform how to analyze the data

- ▶ visualization
- ▶ identify patterns
- ▶ choose appropriate statistical test

2. Presentation of the data

- ▶ summarize results
- ▶ communicate results
- ▶ publish results



Types of graphs - *Exploratory*

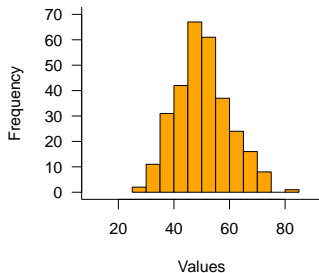
Exploratory graphs help understand the distribution of the data:

- ▶ are the data normally distributed
 - ▶ important assumption in statistics
 - ▶ determines how data are analyzed
- ▶ what is the central tendency
- ▶ what is the spread
- ▶ general summaries of the data

Exploratory: *Histogram*

- ▶ width of bars are defined data bins or intervals
- ▶ height of bars represent bin-specific frequencies

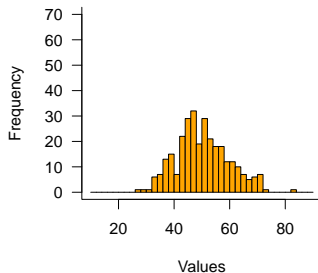
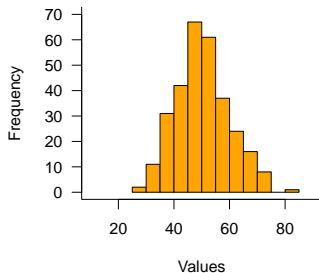
```
hist(values)
```



Exploratory: *Histogram*

- ▶ width of bars are defined data bins or intervals
- ▶ height of bars represent bin-specific frequencies

```
hist(values)  
hist(values, breaks=seq(10,90,2))
```

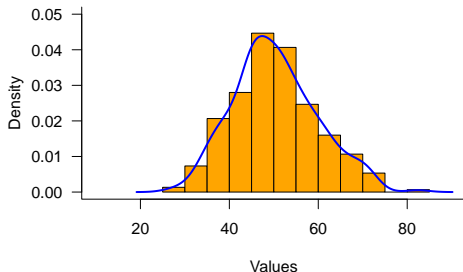


Exploratory: *Histogram + Density Plot*

A *density plot* provides a smooth representation of the histogram:

- ▶ can overlay the density plot
- ▶ requires that a *probability* version of the histogram is plotted

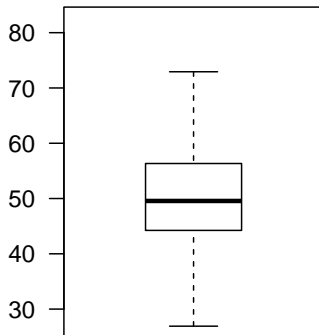
```
hist(values, probability=TRUE)  
lines(density(values))
```



Exploratory: *Box-whisker/Box plot*

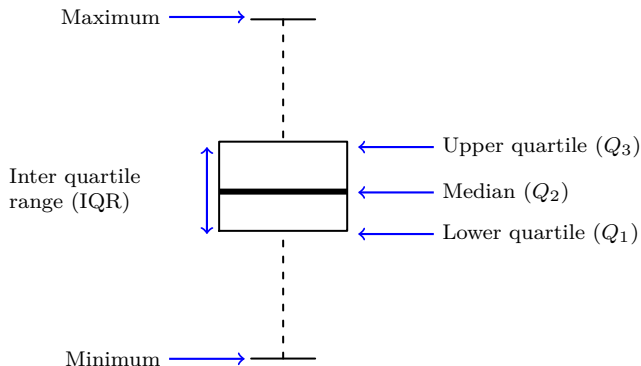
- ▶ distribution
- ▶ outliers
- ▶ symmetry or skewness

```
boxplot(values)
```



Exploratory: *Box-whisker/Box plot*

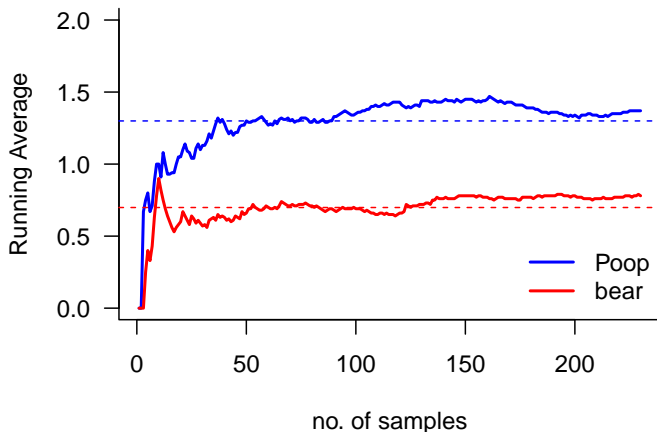
- R: `boxplot(x)` # `x` is data



Exploratory: *Line graph*

Line graph is a useful plot for running average or time series data

```
plot(bear.run, type="l") #"l": line, "p": points, "b": both  
lines(poop.run)
```



Differences

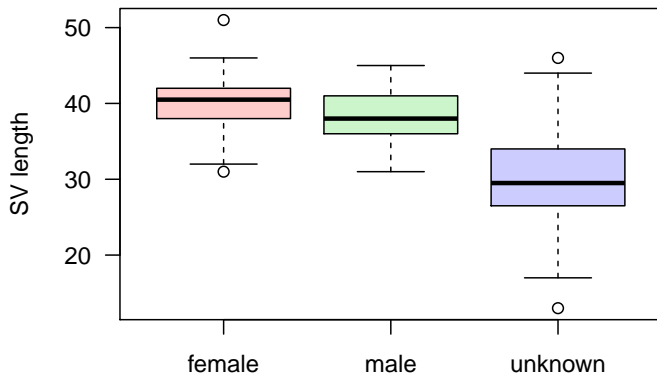
To visualize differences between groups

- ▶ box-whisker plots
 - ▶ compares averages
 - ▶ compares distribution
- ▶ bar charts
 - ▶ compares averages

Differences: *Box-whisker plot*

Compare salamander snout-vent lengths by three sexes:

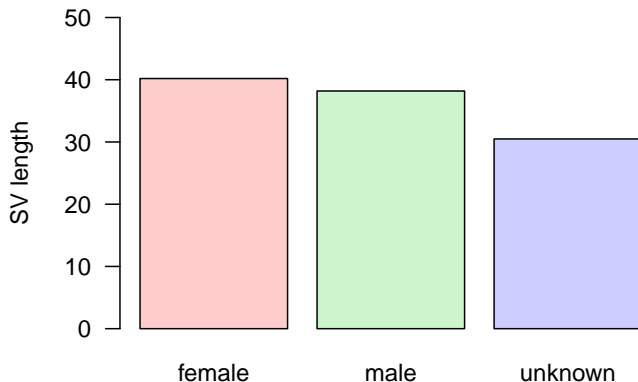
```
boxplot(mander$SVL ~ mander$Sex) #formula notation
```



Differences: *Bar chart*

Compare salamander snout-vent lengths by three sexes:

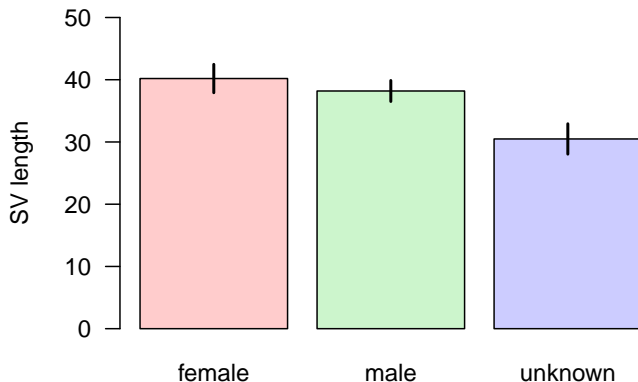
```
bars <- tapply(mander$SVL, mander$Sex, mean) #create matrix (like pivot table)
barplot(bars)                               #plot it
```



Differences: *Bar chart* with associated error

Compare salamander snout-vent lengths by three sexes:

```
bars <- tapply(mander$SVL, mander$Sex, mean)
barplot(bars)
```



Links

Two main approaches for relationships between data:

1. Correlations
2. Associations

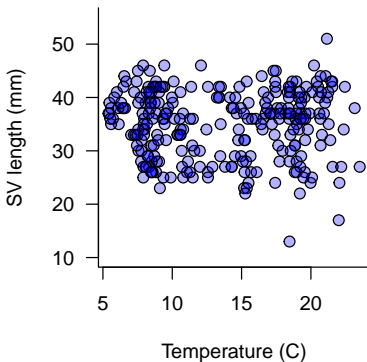
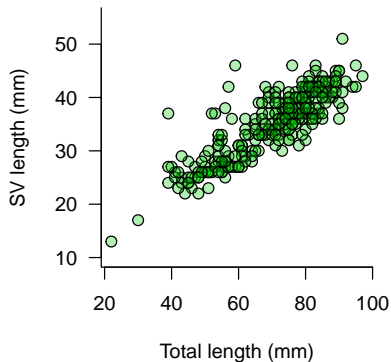
Two main approaches for graphing relationships between data:

1. Correlations

- ▶ two numeric variables
 - ▶ *dependent* variable (of primary interest: y-axis)
 - ▶ *independent* variable (explanatory variable: x-axis)
- ▶ how one variable is related to another
- ▶ *scatter plots*

Links: *Scatter plot*

```
plot(x,y) # x and y are numeric vectors
```



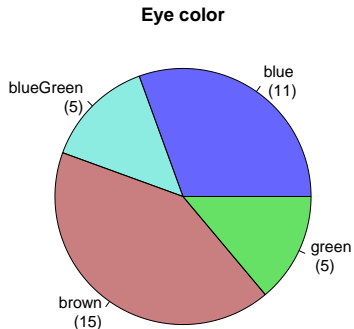
Two main approaches for graphing relationships between data:

2. Associations

- ▶ categorical data
- ▶ summarize categories
 - ▶ counts
 - ▶ proportions
 - ▶ by rows and/or columns of a table
- ▶ *pie charts* for single categories
- ▶ *bar graphs* for several categories

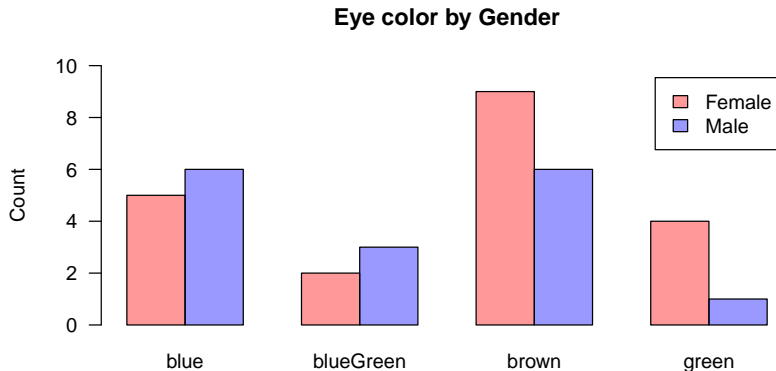
Links: *Pie chart*

```
pietab <- table(classData$Eyes)
pie(pietab) #(number of people with each eye color)
```



Links: *Bar chart*

```
bartab <- table(classData$Gender,classData$Eyes)  
barplot(bartab, beside=TRUE)  #(number of each gender with each eye color)
```



Some graphics pointers

In summary, graphs are a useful data visualization tool

- ▶ summarizing
- ▶ understanding
- ▶ describing
- ▶ presenting/communicating

Some graphics pointers

In summary, graphs are a useful data visualization tool

- ▶ summarizing
- ▶ understanding
- ▶ describing
- ▶ presenting/communicating

BUT we must label the well or they are useless!

- ▶ label both axes
- ▶ provide a main title for your graph
- ▶ avoid clutter
- ▶ make it readable
- ▶ *I expect graphs to be properly labeled from now on!*

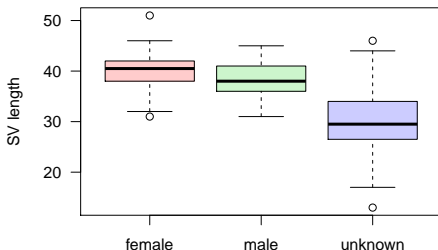
Some graphics pointers

In summary, graphs are a useful data visualization tool

Purpose	Graph Type
Illustrating <i>distribution</i>	Histogram, Density plot Box(-whisker) plot
Illustrating <i>differences</i>	Bar chart, Box plot
Illustrating <i>correlations</i>	Scatter plot
Illustrating <i>associations</i>	Pie chart, Bar chart
Illustrating <i>sample size</i>	Line plot of running avg

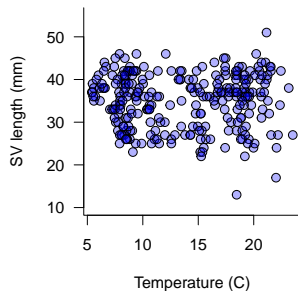
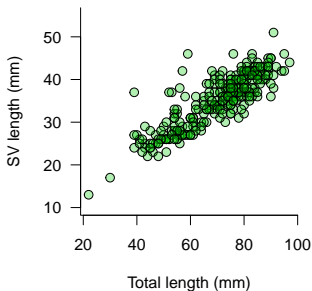
Beyond graphs, Towards statistics

- ▶ Graphs are powerful tools that provide insight and understanding of the patterns and relationships in the data.
- ▶ Don't give us the answer though:
 - ▶ are differences *significant*?
 - ▶ are associations *significant*?



Beyond graphs, Towards statistics

- ▶ Graphs are powerful tools that provide insight and understanding of the patterns and relationships in the data.
- ▶ Don't give us the answer though:
 - ▶ are differences *significant*?
 - ▶ are associations *significant*?



Beyond graphs, Towards statistics

- ▶ Graphs are powerful tools that provide insight and understanding of the patterns and relationships in the data.
- ▶ Don't give us the answer though:
 - ▶ are differences *significant*?
 - ▶ are associations *significant*?
- ▶ Statistics is the tool we use to formally answer these questions!
 - ▶ the differences *are/are not* significant!
 - ▶ are associations *are/are not* significant!

Group evaluations

1. Are you sitting with your group?

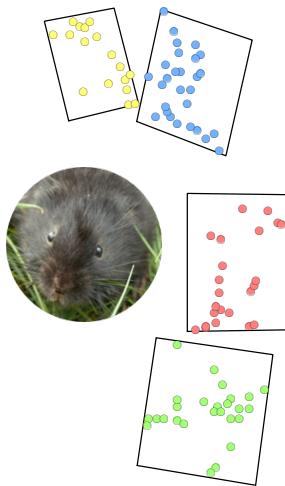
A) Yes

B) No

Group practical: water vole weights

Ultimately we are interested in comparing sex-specific water vole weights across multiple populations (networks). The data include weight measurements of:

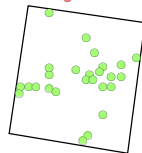
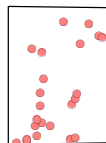
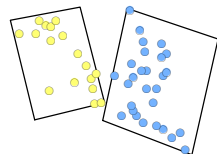
- ▶ 100's voles
- ▶ from 4 water vole sub-populations
- ▶ from males and females



Group practical: water vole weights

Ultimately we are interested in comparing sex-specific water vole weights across multiple populations (networks). The data include weight measurements of:

- ▶ 100's voles
- ▶ from 4 water vole sub-populations
- ▶ from males and females



The assignment:

- ▶ download data & empty script
- ▶ complete the script (in groups)
- ▶ submit to moodle (1 per group)