

Graphs:
The art of designing information

"A picture tells a thousand words"

- Lake Blanche

1

Graphs are used to try to tell a story



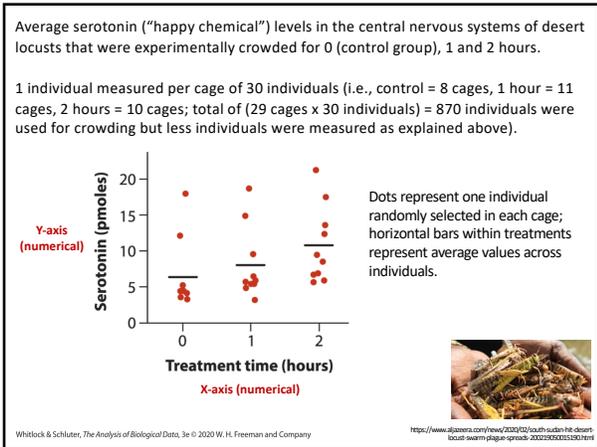
...and to make a point

2

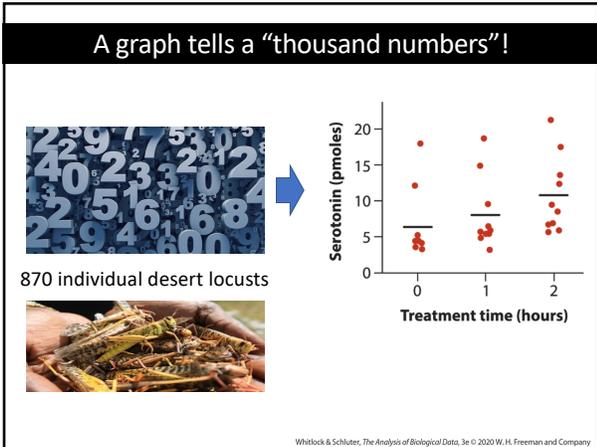
General definition of a graph

- Visual representation of a relationship between two or three variables (and more sometimes).
- Variables can be of any type (e.g., categorical or numerical).
- They commonly consist of two axes: x-axis (horizontal or abscissa) and y-axis (vertical or ordinate).

3



4



5

Why graphs?

- Powerful way of summarizing data that is easy to read (i.e., quick and direct).
- Highlight the most important information (i.e., facilitate communication).
- Help understand the data.
- Reveal structure and patterns in the data
- Help convince others.
- Easy to remember (general trends).
- Aid in detecting unusual features in data.
- Tell stories.

6

Types of graphs

There are lots of types of graphs! The most commons (and covered in BIOL322) are:

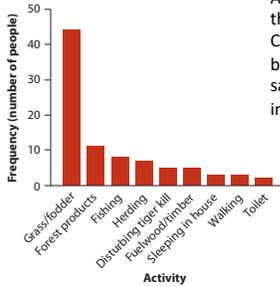
TODAY:

- Bar graph
- Pie chart
- Histogram
- Line graph
- Scatter plot
- Strip chart
- Graphs of data distributions (box plots, histograms, violin plot)



7

BAR GRAPH: Vertical or horizontal columns (bars) representing the distribution of a numerical variable against one or more categoribal variable.



Activities of people at the time they were killed by tigers near Chitwan National Park (Nepal) between 1979-2006; n=88 (n = sample size, i.e., number of killed individuals).

Activity - categorical
Frequency - numerical (discrete)

Whitlock & Schluter, The Analysis of Biological Data, 3e © 2020 W. H. Freeman and Company

8

BAR GRAPHS are usually better than pie charts

Activities of people at the time they were killed by tigers near Chitwan National Park (Nepal) between 1979-2006; n=88



Whitlock & Schluter, The Analysis of Biological Data, 3e © 2020 W. H. Freeman and Company

9

BAR GRAPH: Two categorical variables
(often from a contingency table)

Is reproduction associated with health risks?

	Control group	Egg-removal group	Row total
Malaria	7	15	22
No Malaria	28	15	43
Column total	35	30	65

Parus major Treatment (egg removal/control) & outcome (malaria – yes/no) - categorical
Frequency - numerical (discrete)



Female birds put more energy in generating eggs to make up for those removed, thus reducing energy allocation towards immunocompetence.

10

BAR GRAPH: Two categorical variables
(often from a contingency table)

Is reproduction (explanatory variable = egg removal) associated with health risks (response variable = malaria susceptibility)?

	explanatory variable		Row total
	Control group	Egg-removal group	
Malaria	7	15	22
No Malaria	28	15	43
Column total	35	30	65

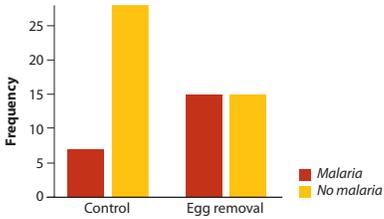
Parus major Treatment (egg removal/control) & outcome (malaria – yes/no) - categorical
Frequency - numerical (discrete)



egg removal forces females to produce additional eggs (i.e., increase reproduction)

11

BAR GRAPH: Two categorical variables
(often from a contingency table)



Treatment (egg removal/control) & outcome (malaria – yes/no) - categorical
Frequency - numerical (discrete)

12

BAR GRAPH: Two categorical variables
(often from a contingency table)

Treatment	Malaria	No malaria
Control	7	27
Egg removal	15	15

For bar graphs, it's generally recommended to start the measurement axis at zero to ensure the relative sizes of the bars accurately reflect the data (more on this issue at the end of this lecture).

13

This example illustrates different study designs and how biologists infer cause and effect.

Treatment	Malaria	No malaria
Control	7	27
Egg removal	15	15

Treatment (egg removal/control) & outcome (malaria – yes/no) - categorical
Frequency - numerical (discrete)

14

Explanatory versus Response variables

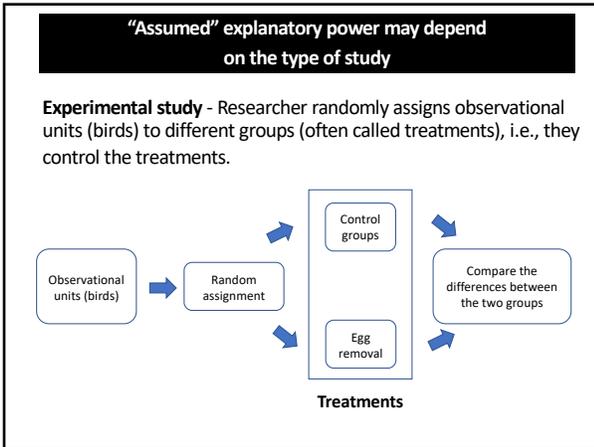
A major goal of Biostatistics is to relate one variable to another, by examining associations between variables or differences among groups.

When studying associations, we often ask how well one variable - the explanatory variable—helps predict or explain another variable - the response variable.

Is reproduction (explanatory variable = egg removal) associated with health risks (response variable = malaria susceptibility)?

“Assumed” explanatory power may depend on the type of study:
[1] experimental versus [2] observational studies

15



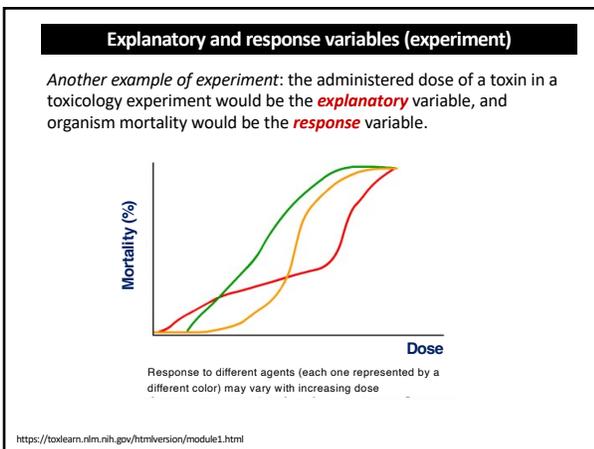
16

Explanatory and response variables (experiment)

When conducting an *experiment* (e.g., malaria study in the last slides), the treatment variable (the one manipulated by the researcher, i.e., egg removal) is the **explanatory** variable, and the measured effect of the treatment (malaria susceptibility) is the **response** variable.

		explanatory variable		Row total
		Control group	Egg-removal group	
response variable	Malaria	7	15	22
	No Malaria	28	15	43
	Column total	35	30	65

17



18

"Assumed" explanatory power may depend on the type of study

Observational study - Researchers have no control over which observational units fall into which treatment or values of the explanatory variable. Examples:

- Studies on the health consequences of cigarette smoking in humans (unethical to assign smoking and no-smoking treatments to observational units, i.e., people).
- Growth of fish in warm versus cold lakes (observational units, i.e., fish are already in lakes; the research has no control on which fish goes in which lake).

19

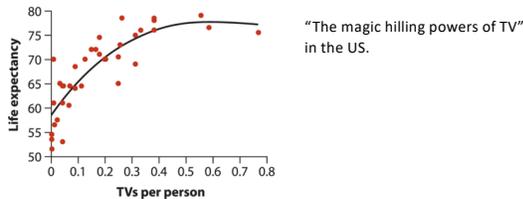
Let's take a break – 1 minute



20

Explanatory and response variables (observational study)

When neither variable is manipulated by the researcher (i.e., observational study; sample of convenience), their association might nevertheless be described by the "effect" of one of the variables (the explanatory) on the other (the response), even though the association itself is not direct evidence for causation.



Whitlock & Schluter, *The Analysis of Biological Data*, 3e © 2020 W. H. Freeman and Company

21

Explanatory and response variables (observational study)

When neither variable is manipulated by the researcher (i.e., observational study; sample of convenience), their association might nevertheless be described by the "effect" of one of the variables (the explanatory) on the other (the response), even though the association itself is not direct evidence for causation.

"The magic hilling powers of TV" in the US

Overall wealth of citizens through time (and cheaper TVs)

"cause" "cause"

TVs/person correlation Life expectancy

The "causal" mechanism here could be many things, including better access to health care.

Whitlock & Schluter, The Analysis of Biological Data, 3e © 2008 W. H. Freeman and Company

22

Explanatory and response variables (observational study - common trends may not mean cause & effect, i.e., they are independent variables)

Divorce rate in Maine correlates with Per capita consumption of margarine (US)

	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009
Divorce rate in Maine (Divorces per 1000 people (US Census))	5	4.7	4.6	4.4	4.3	4.1	4.2	4.2	4.2	4.1
Per capita consumption of margarine (US) (Pounds (USDA))	8.2	7	6.5	5.3	5.2	4	4.6	4.5	4.2	3.7

Correlation: 0.992558

[Permalink - Mark as interesting - Not interesting](http://davidgen.com/view_correlation%3d123)

23

Independent versus dependent variables = explanatory versus response variables, respectively

Strictly speaking, if one variable depends on another, neither is truly "independent"; instead, we refer to them as the explanatory and response variables.

Sometimes you will hear variables referred to as "**independent**" and "**dependent**". These are the same as **explanatory** and **response** variables, respectively.

24

Independent versus dependent variables = Explanatory versus response variables, respectively

Regardless whether the association is causal, the "expected" explanatory variable goes in the X-axis and the expected response variable goes in the Y-axis.

The scatter plot shows a positive correlation between the number of TVs per person (X-axis, 0 to 0.8) and life expectancy (Y-axis, 50 to 80). A smooth curve is fitted to the data points, showing that life expectancy increases with the number of TVs but at a decreasing rate.

Whitlock & Schluter, The Analysis of Biological Data, 3e © 2020 W. H. Freeman and Company

25

Back to BAR GRAPHS: two categorical variables

Is reproduction associated with health risks?
Not so clear from this bar graph

The bar graph shows the frequency of malaria (red bars) and no malaria (yellow bars) for two treatments: Control and Egg removal. For the Control group, there are approximately 7 malaria cases and 27 no malaria cases. For the Egg removal group, there are approximately 15 malaria cases and 15 no malaria cases. Blue arrows point to the bars, and an equals sign is placed between the two bars in the Egg removal group.

Treatment (egg removal/control) & outcome (malaria – yes/no) - categorical
Frequency - numerical (discrete)

26

BAR GRAPHS (stacked = mosaic graph): Two categorical variables

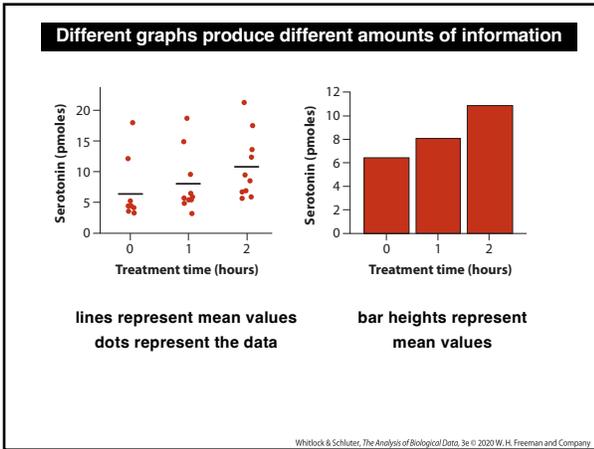
Is reproduction risky to health? Much clearer now!

The stacked bar graph shows the relative frequency of malaria (red) and no malaria (yellow) for two treatments: Control and Egg removal. For the Control group, the relative frequency of malaria is approximately 0.25 and no malaria is 0.75. For the Egg removal group, the relative frequency of malaria is approximately 0.5 and no malaria is 0.5.

Treatment (egg removal/control) & outcome (malaria – yes/no) - categorical
Frequency - numerical (discrete)

Whitlock & Schluter, The Analysis of Biological Data, 3e © 2020 W. H. Freeman and Company

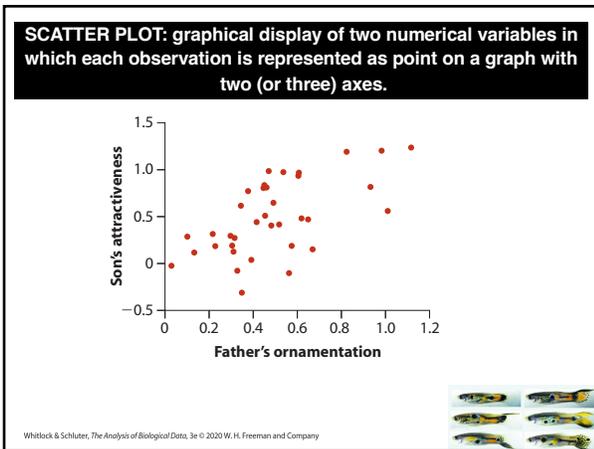
27



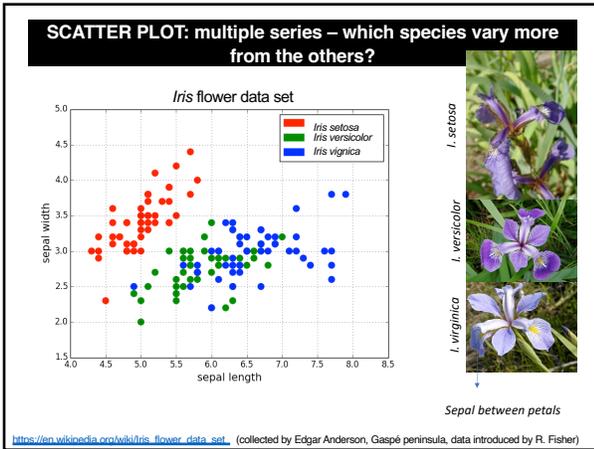
31



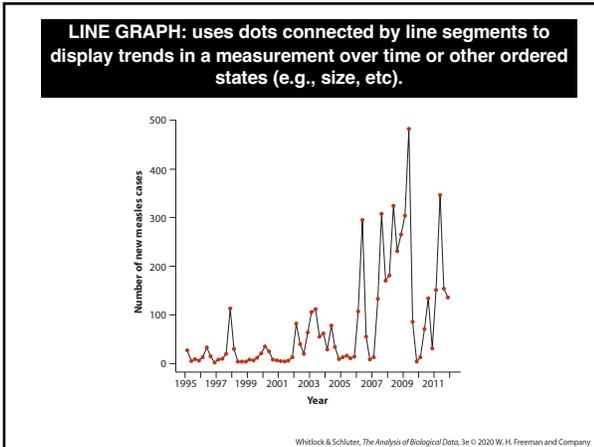
32



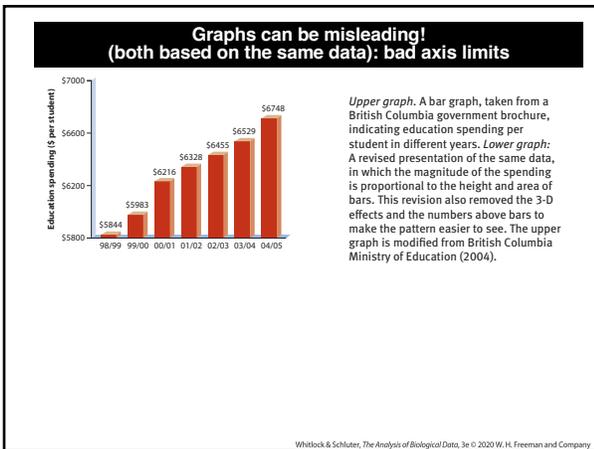
33



34



35



36

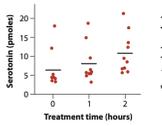
Mistake 1: Hide the data

How to hide data:

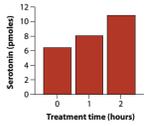
- Provide only statistical summaries (e.g., means).

How to reveal data:

- Present all data points, while allowing all to be seen.



lines represent mean values
dots represent the data



bar heights represent mean values

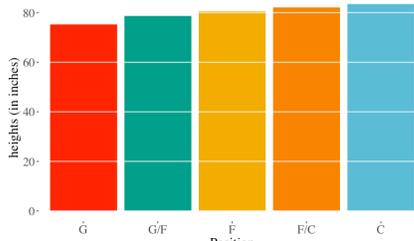
© 2020 W.H. Freeman and Company

40

Not Showing Data, Just Summaries

This plot hides the variation within positions (only means).

Mean heights of NBA players by position



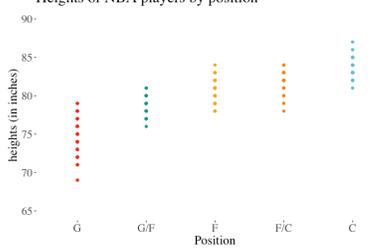
center, forward, guard © 2020 W.H. Freeman and Company

41

Not Showing Data, Over-Plotting

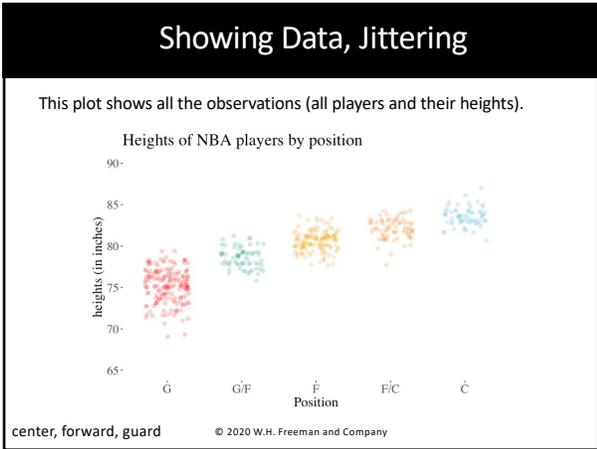
This plot hides the density of observations (number of players, i.e., many have very similar heights).

Heights of NBA players by position



center, forward, guard © 2020 W.H. Freeman and Company

42



43

Mistake 2: Making Patterns Hard to See

How to hide patterns:

- Make one plot and call it good.
- Use unreasonable scales.
- Arrange factors nonsensically.

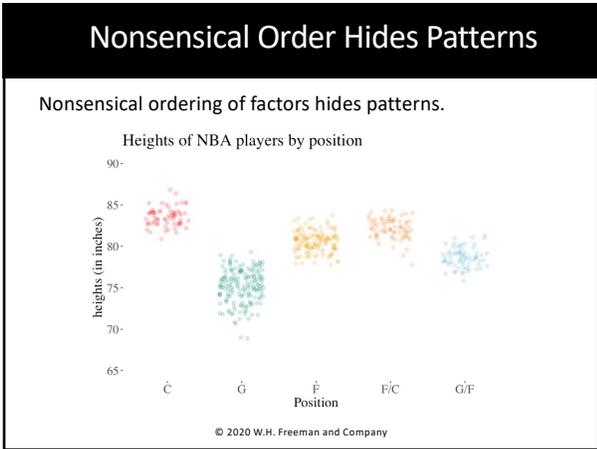
How to reveal patterns:

- Explore multiple potential plots.
- Use appropriate scales.
- Arrange factors meaningfully.

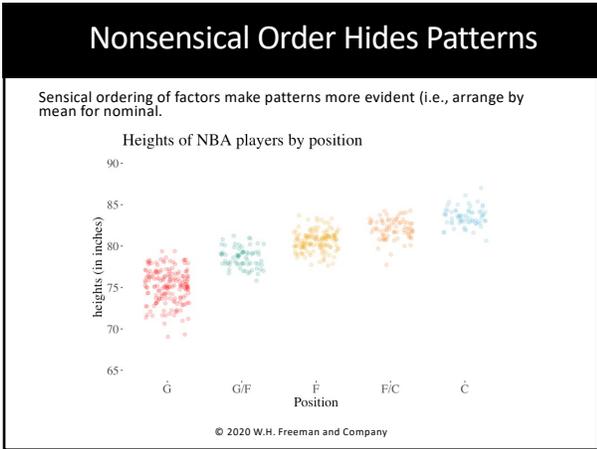
Arrange in order for ordinal, and by mean for nominal.

© 2020 W.H. Freeman and Company

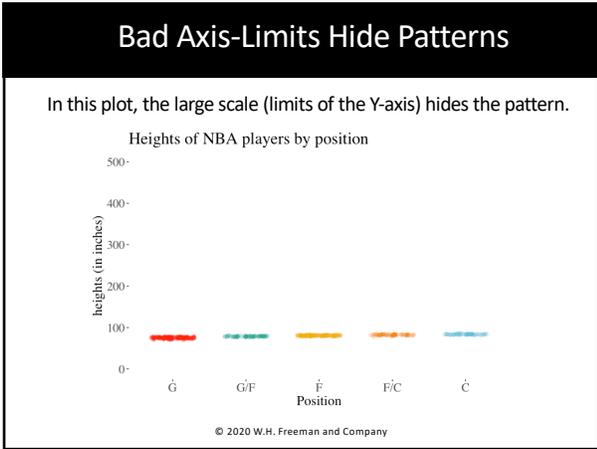
44



45



46



47

Graphs:

The art of designing information

“A picture tells a thousand words”

- Lake Blanche

48

Next lecture: How to build frequency distributions and introduction to descriptive (or summary) statistics

