

Graphs: The art of designing information

“A picture tells a thousand words”

- *Lake Blanche*

Graphs are used to try to tell a story

HERMAN[®]

by Jim Unger



“That’s the last time I go on vacation”

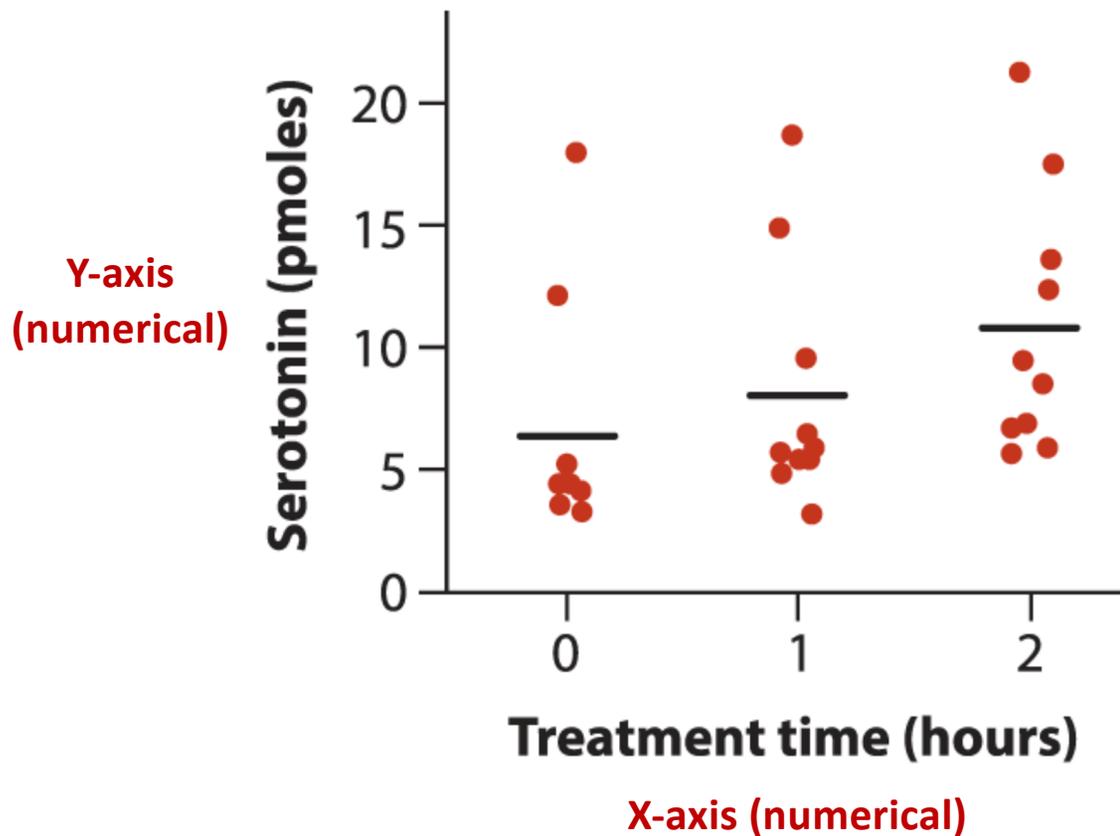
...and to make a point

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).

Average serotonin (“happy chemical”) levels in the central nervous systems of desert locusts that were experimentally crowded for 0 (control group), 1 and 2 hours.

1 individual measured per cage of 30 individuals (i.e., control = 8 cages, 1 hour = 11 cages, 2 hours = 10 cages; total of (29 cages x 30 individuals) = 870 individuals were used for crowding but less individuals were measured as explained above).



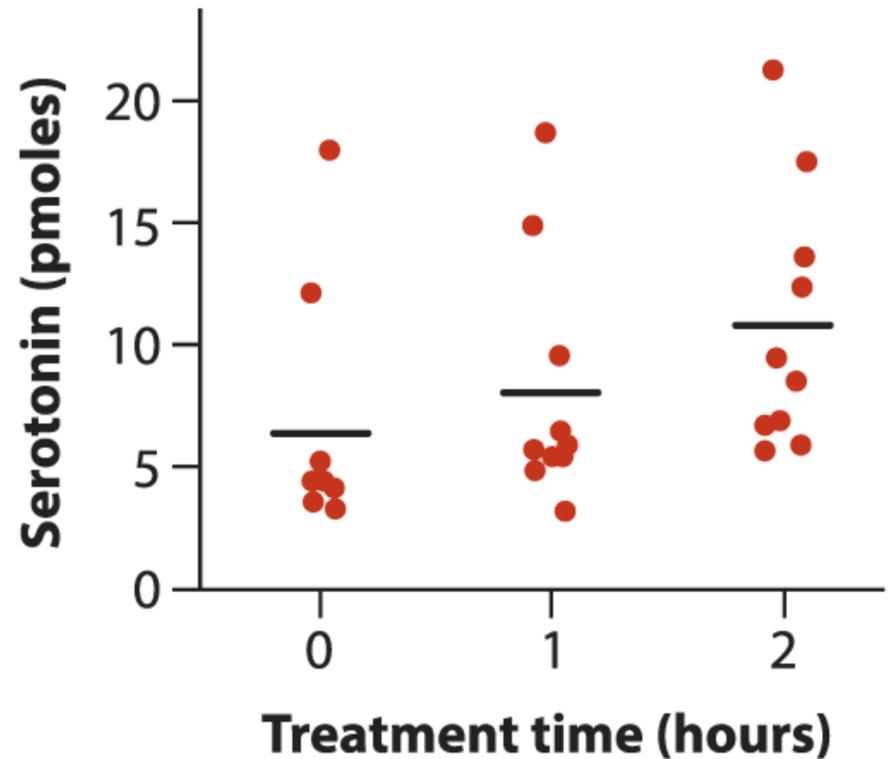
Dots represent one individual randomly selected in each cage; horizontal bars within treatments represent average values across individuals.



A graph tells a “thousand numbers”!



870 individual desert locusts



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.**

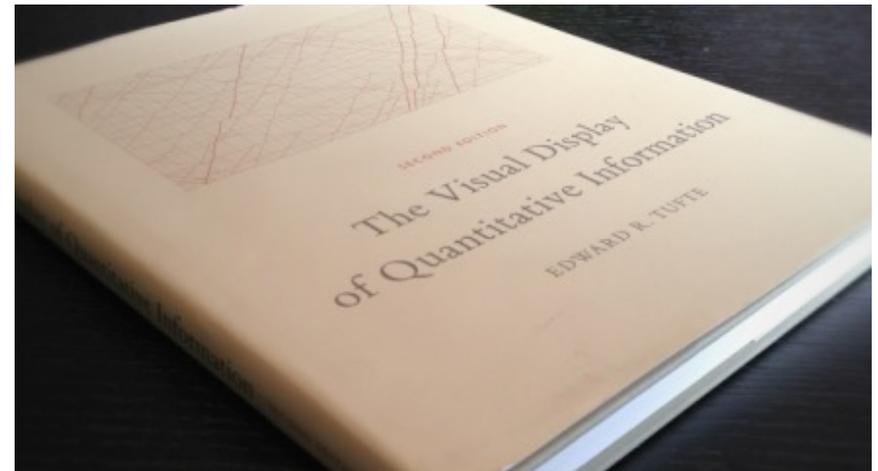


Types of graphs

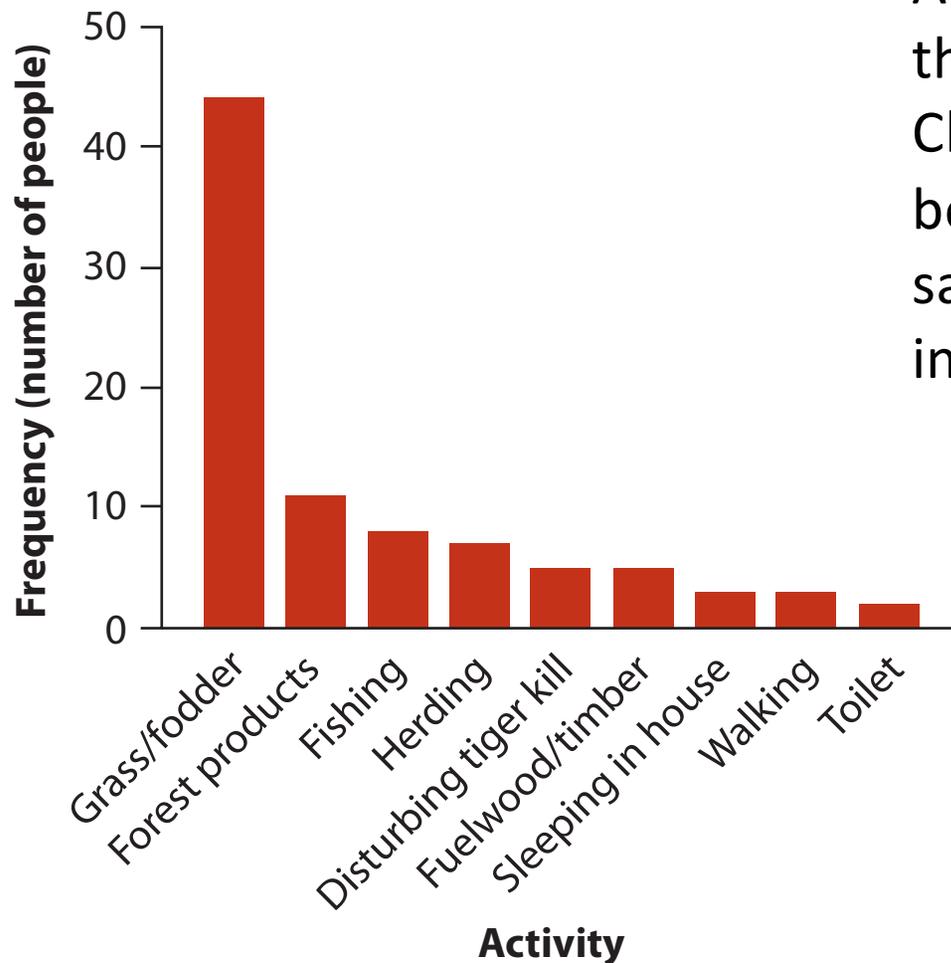
There are lots of types of graphs! The most commons (and covered in BIOL322) are:

TODAY:

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



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

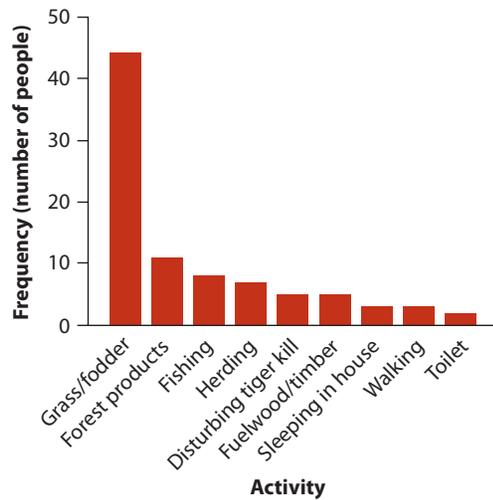


Activities of people at the time they were killed by tigers near Chitwan National Park (Nepal) between 1979-2006; $n=88$ (n = sample size, i.e., number of killed individuals).

Activity - categorical
Frequency - numerical
(discrete)

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



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

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

Frequency - numerical (discrete)

Parus major



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

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	
response variable	Malaria	7	15	22
	No Malaria	28	15	43
	Column total	35	30	65

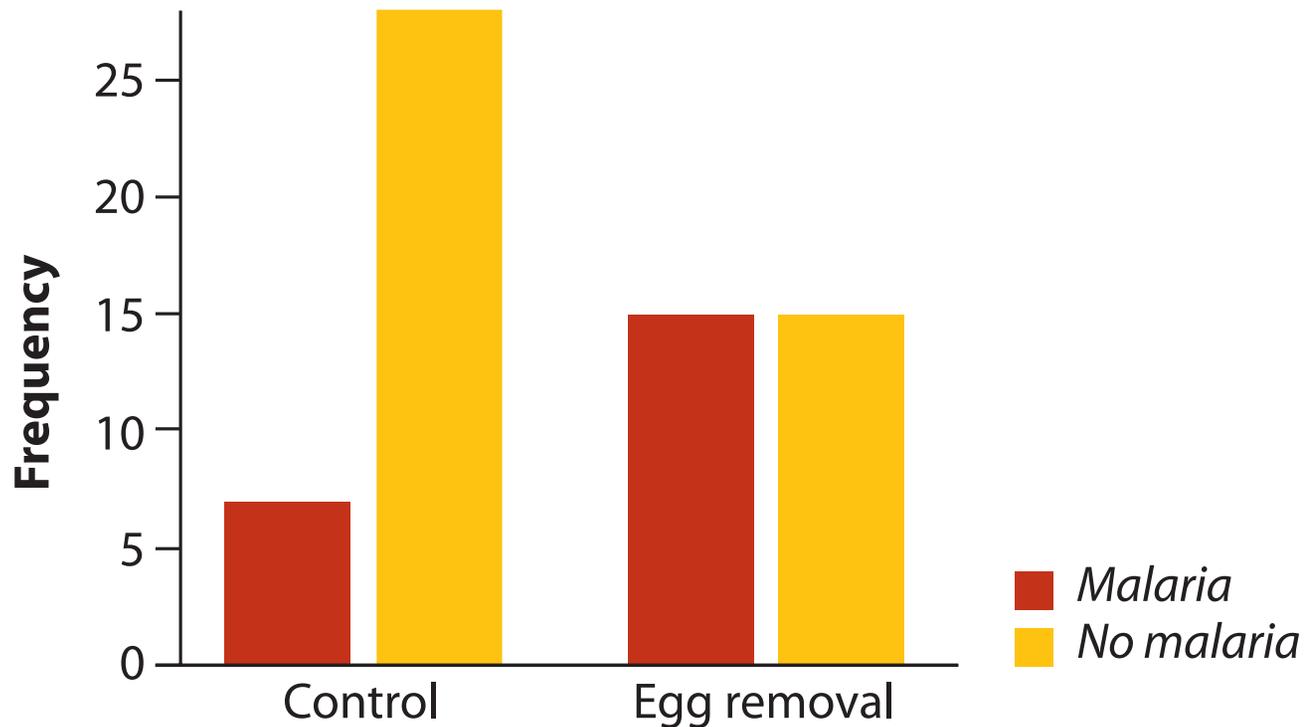
Parus major



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

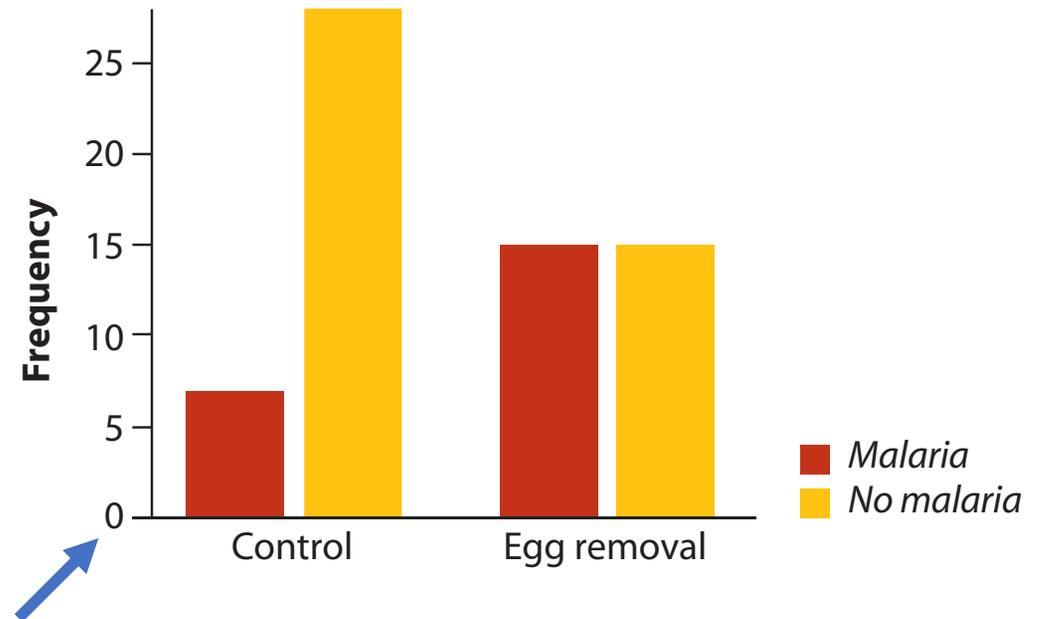
egg removal forces females to produce additional eggs
(i.e., increase reproduction)

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



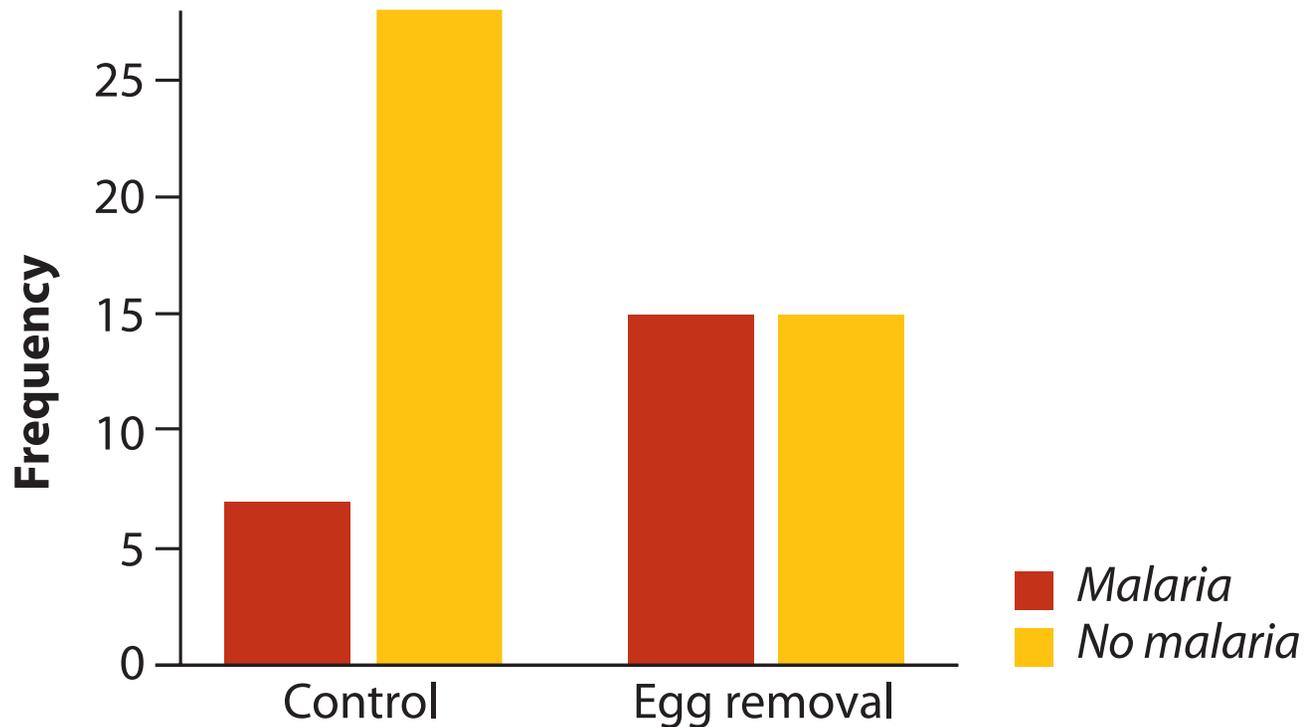
Treatment (egg removal/control) & outcome (malaria – yes/no) - categorical
Frequency - numerical (discrete)

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



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).

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



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

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.

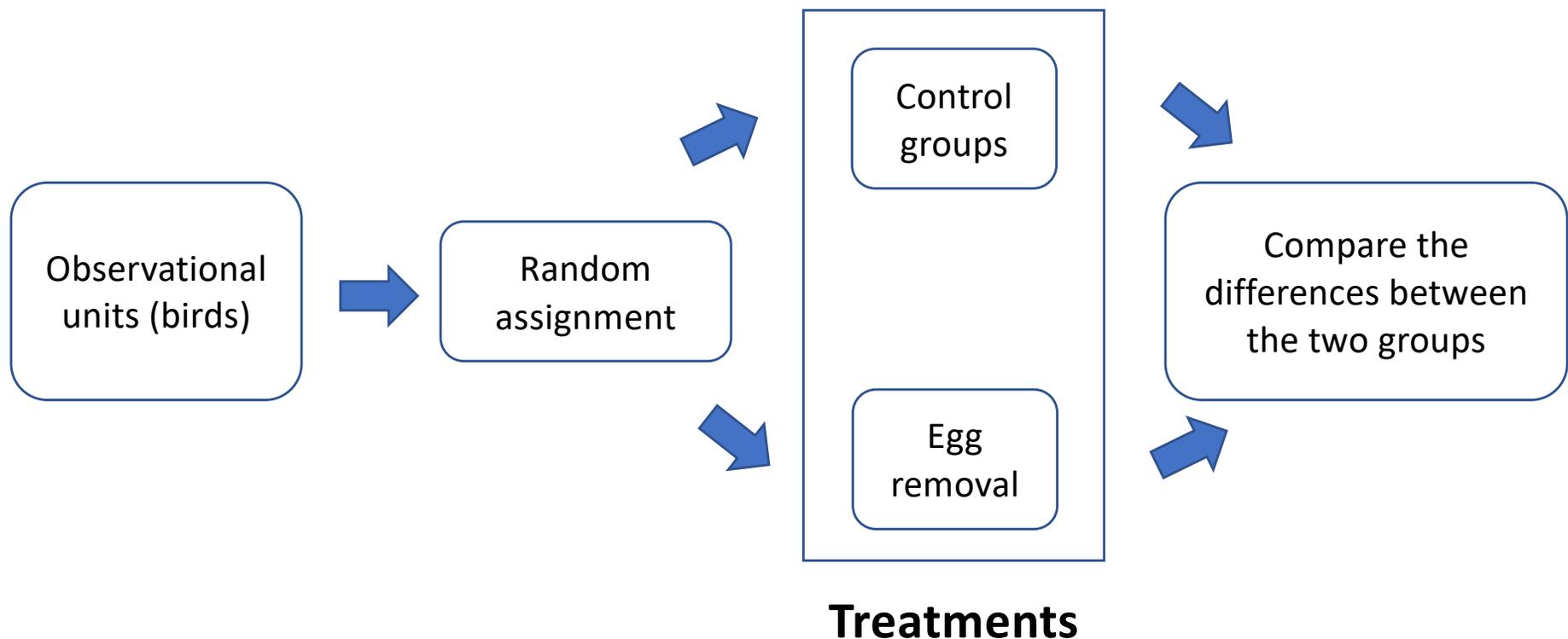
Is reproduction (explanatory variable = egg removal) associated with health risks (response variable = malaria susceptibility)?

“Assumed” explanatory power may depend on the type of study:

[1] **experimental** versus [2] **observational** studies

“Assumed” explanatory power may depend on the type of study

Experimental study - Researcher randomly assigns observational units (birds) to different groups (often called treatments), i.e., they control the treatments.



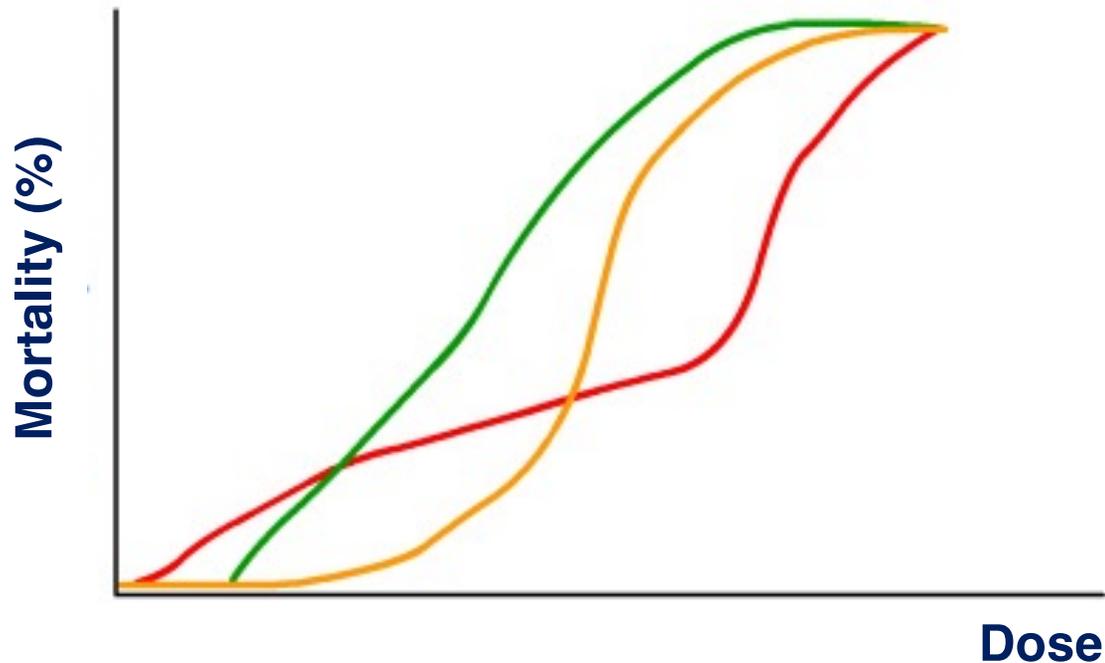
Explanatory and response variables (experiment)

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

Explanatory and response variables (experiment)

Another example of experiment: the administered dose of a toxin in a toxicology experiment would be the **explanatory** variable, and organism mortality would be the **response** variable.



Response to different agents (each one represented by a different color) may vary with increasing dose

“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:

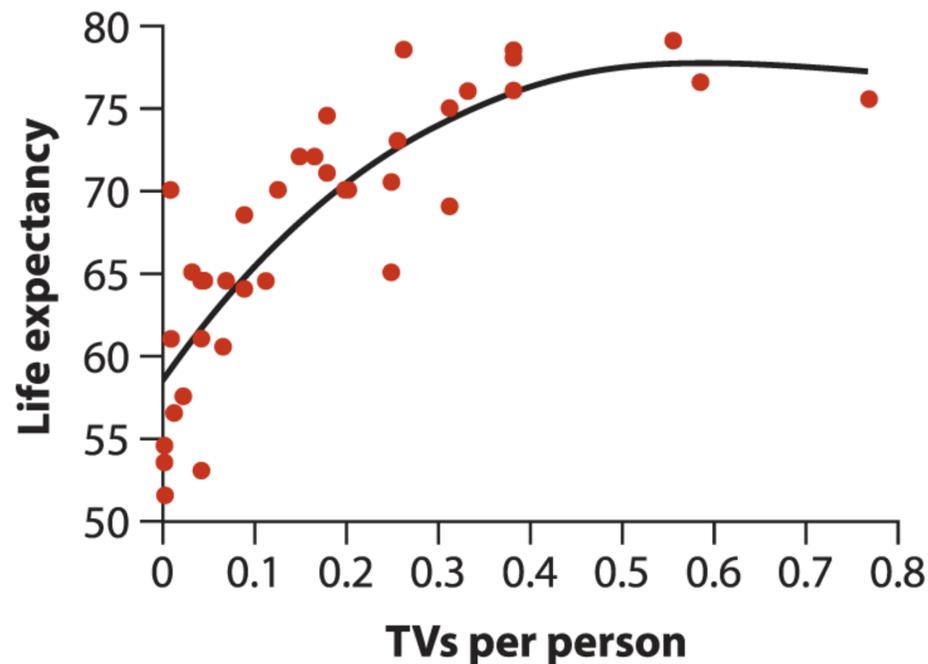
- Studies on the health consequences of cigarette smoking in humans (unethical to assign smoking and no-smoking treatments to observational units, i.e., people).
- Growth of fish in warm versus cold lakes (observational units, i.e., fish are already in lakes; the research has no control on which fish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

Let's take a break – 1 minute



Explanatory and response variables (observational study)

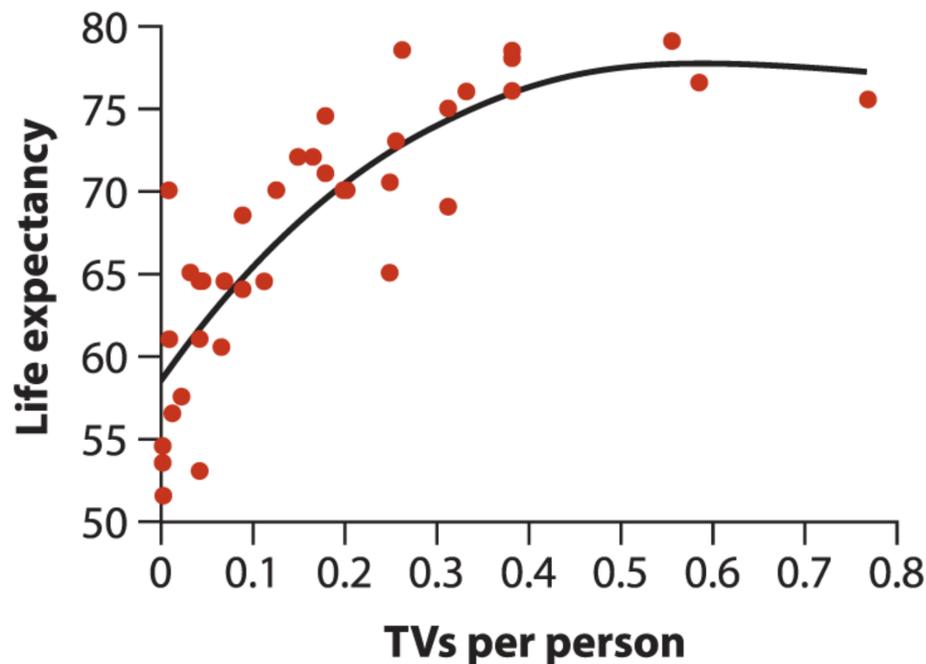
When neither variable is manipulated by the researcher (i.e., observational study; sample of convenience), their association might nevertheless be described by the “effect” of one of the variables (the explanatory) on the other (the response), even though the association itself is not direct evidence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



“The magic hilling powers of TV”
in the US.

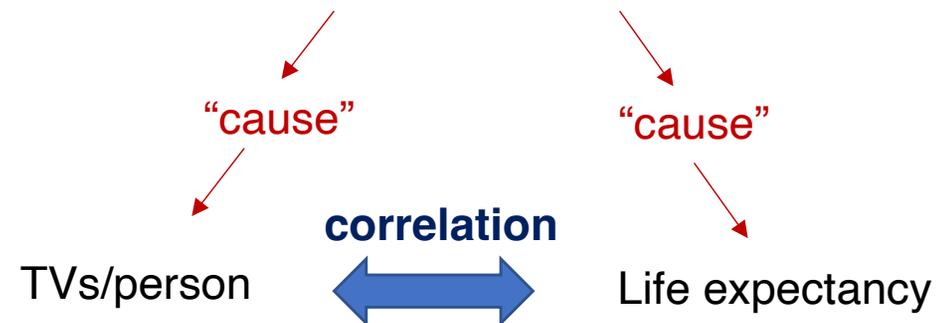
Explanatory and response variables (observational study)

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

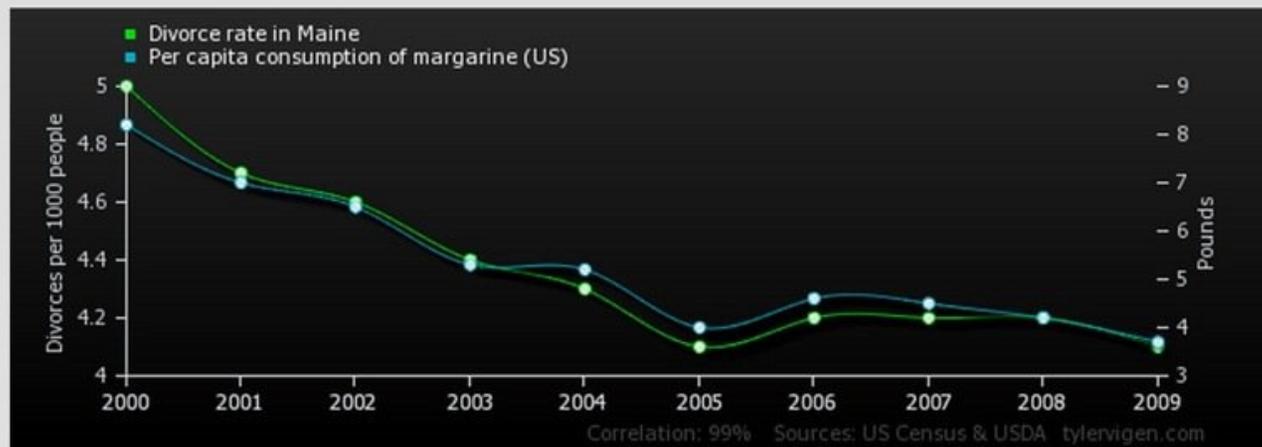
Overall wealth of citizens
through time (and cheaper TVs)



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

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
<i>Divorce rate in Maine</i> 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
<i>Per capita consumption of margarine (US)</i> 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](#)

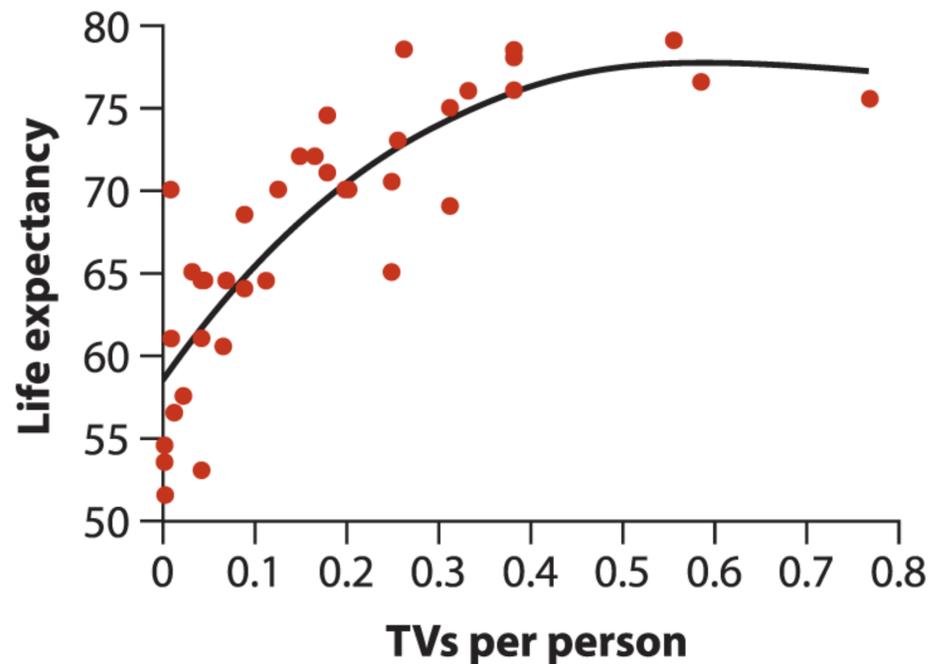
Independent versus dependent variables = explanatory versus response variables, respectively

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.

Independent versus dependent variables = Explanatory versus response variables, respectively

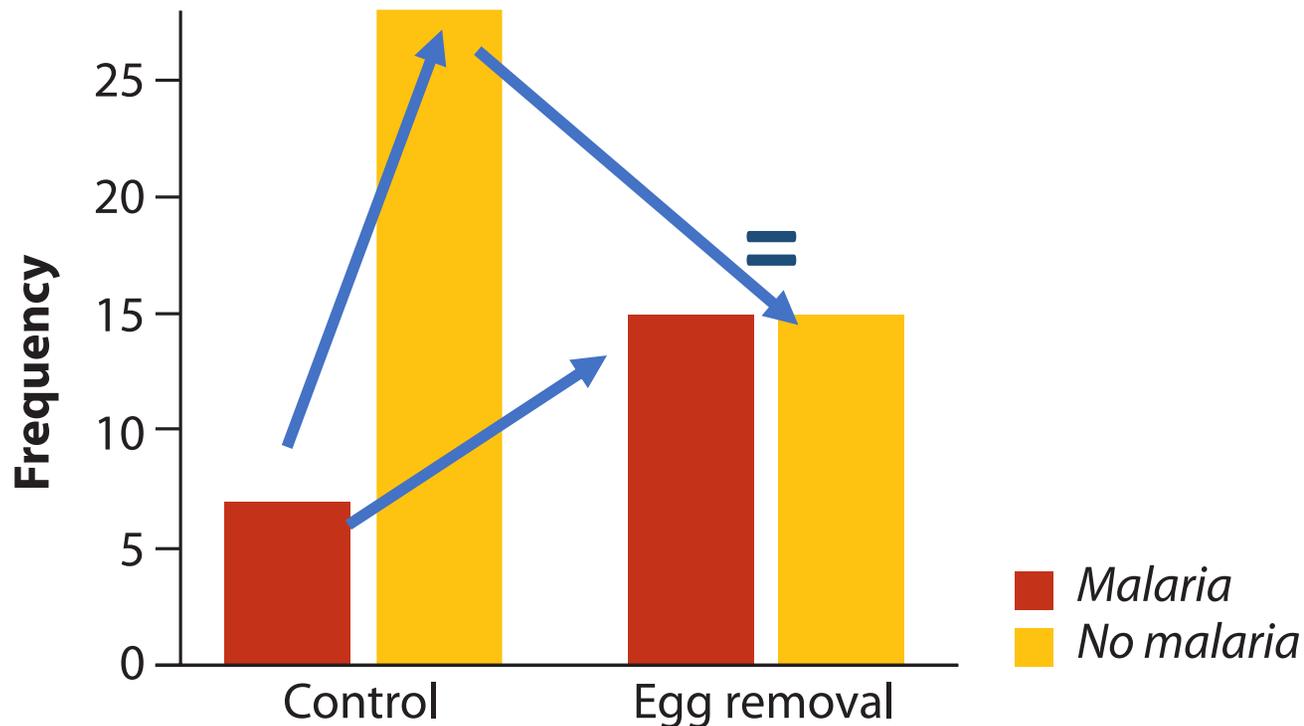
Regardless whether the association is causal, the “expected” explanatory variable goes in the X-axis and the expected response variable goes in the Y-axis.



Back to BAR GRAPHS: two categorical variables

Is reproduction associated with health risks?

Not so clear from this bar graph



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

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

Is reproduction risky to health? Much clearer now!



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

The same dataset can produce multiple visualizations, and each one highlights different patterns.

Grain variety	Cultivation method	Grain yield (t/ha)
A	Broad bed	1.2
A	Traditional	0.7
B	Broad bed	1.6
B	Traditional	1.2
C	Broad bed	1.4
C	Traditional	1.3

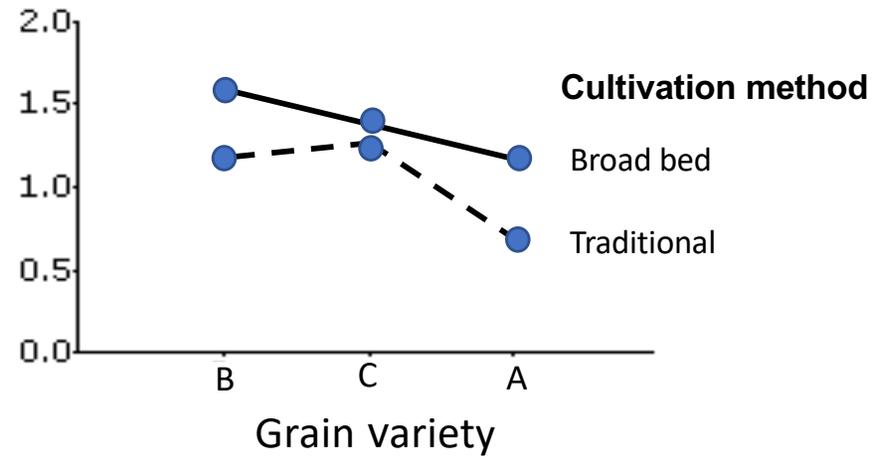
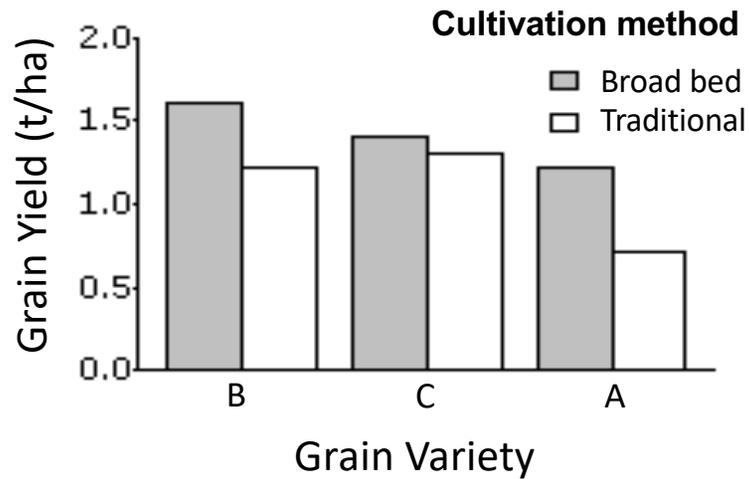


Traditional (continuous; non-spaced)



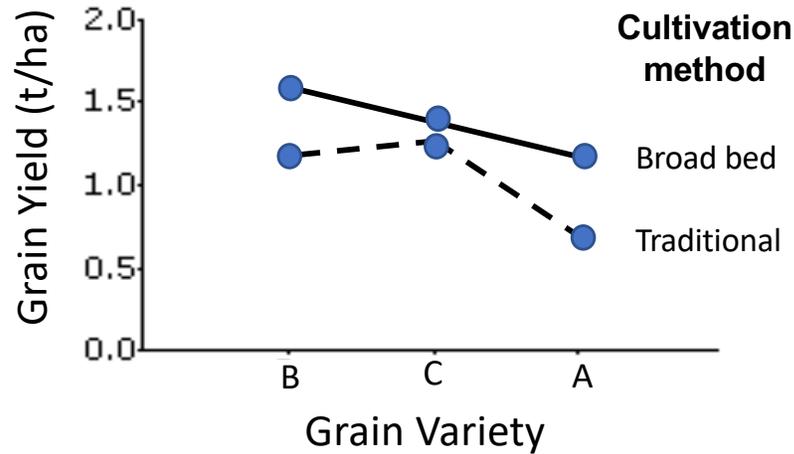
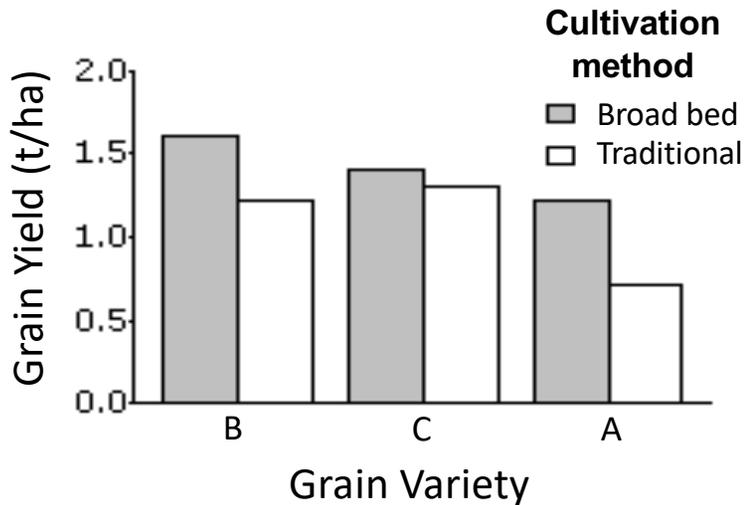
Broad bed (spaced)

BAR GRAPHS are not always the best way! (these graphs are based on the same data)

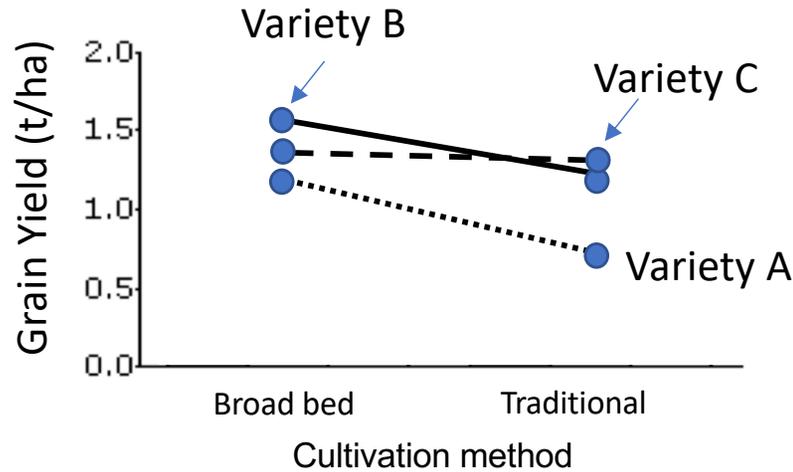


BAR GRAPHS are not always the best way!

which variable goes where? What message do you want to “sell”?



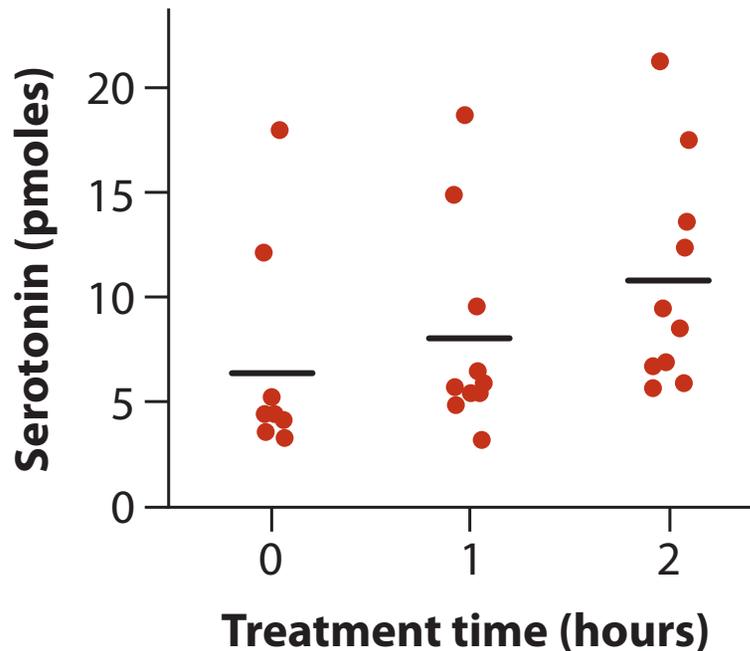
demonstrates better differences between cultivation methods



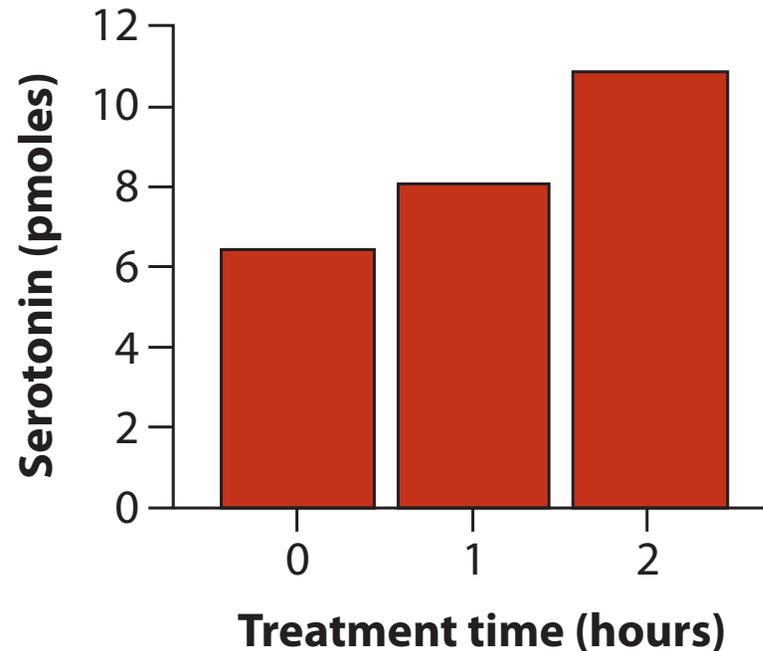
demonstrates better differences between grain varieties

Same exact data in ALL graphs.

Different graphs produce different amounts of information



lines represent mean values
dots represent the data

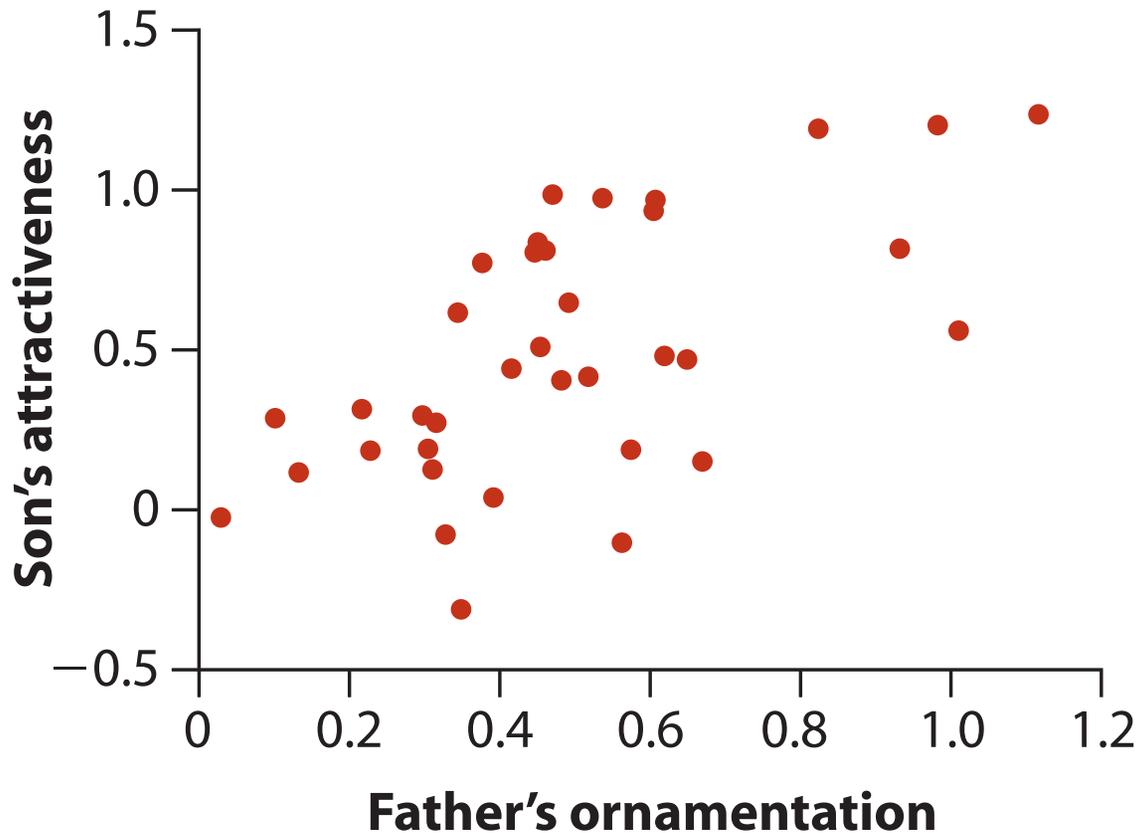


bar heights represent
mean values

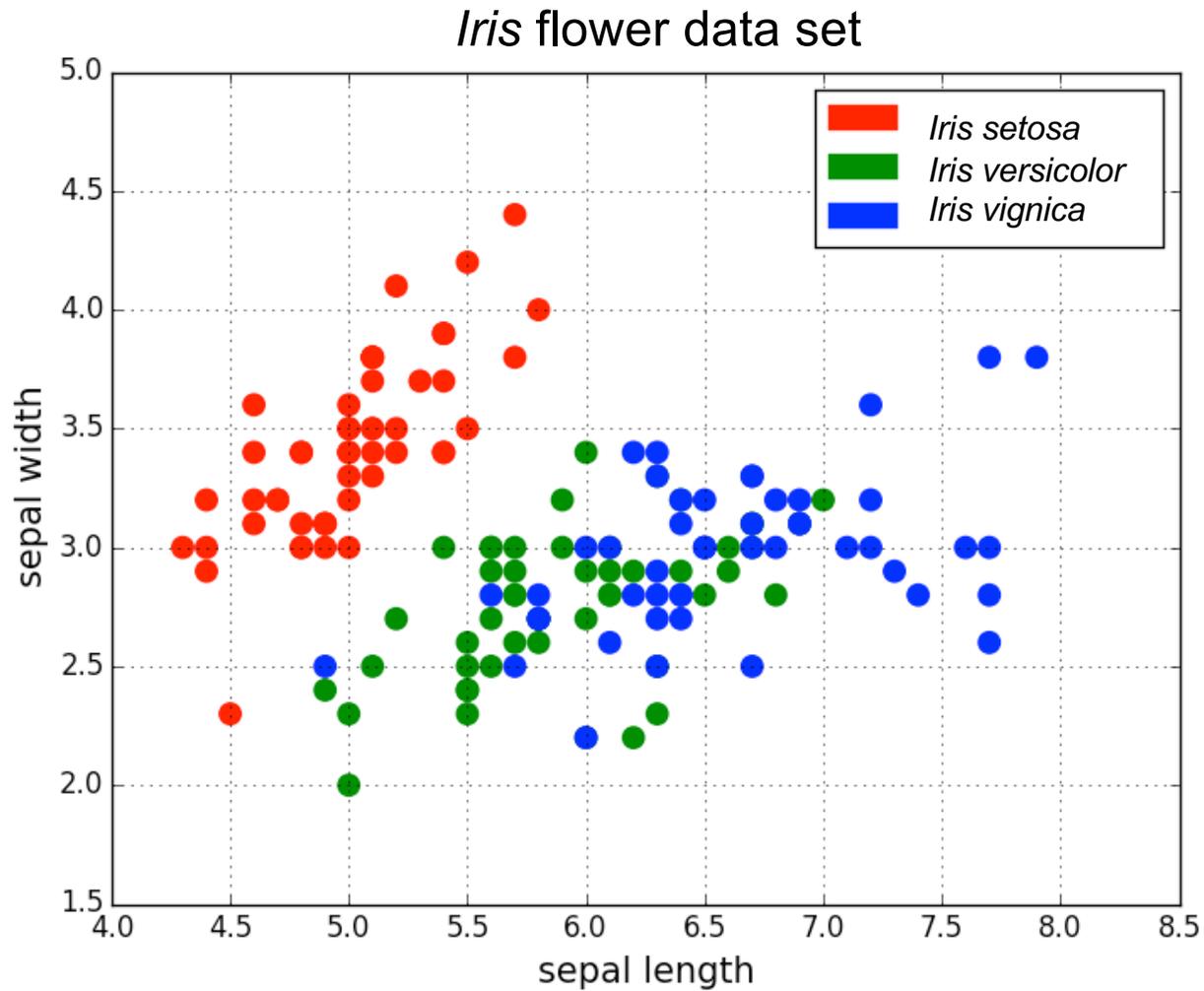
Let's take a break - 1 minute



SCATTER PLOT: graphical display of two numerical variables in which each observation is represented as point on a graph with two (or three) axes.



SCATTER PLOT: multiple series – which species vary more from the others?



I. setosa



I. versicolor

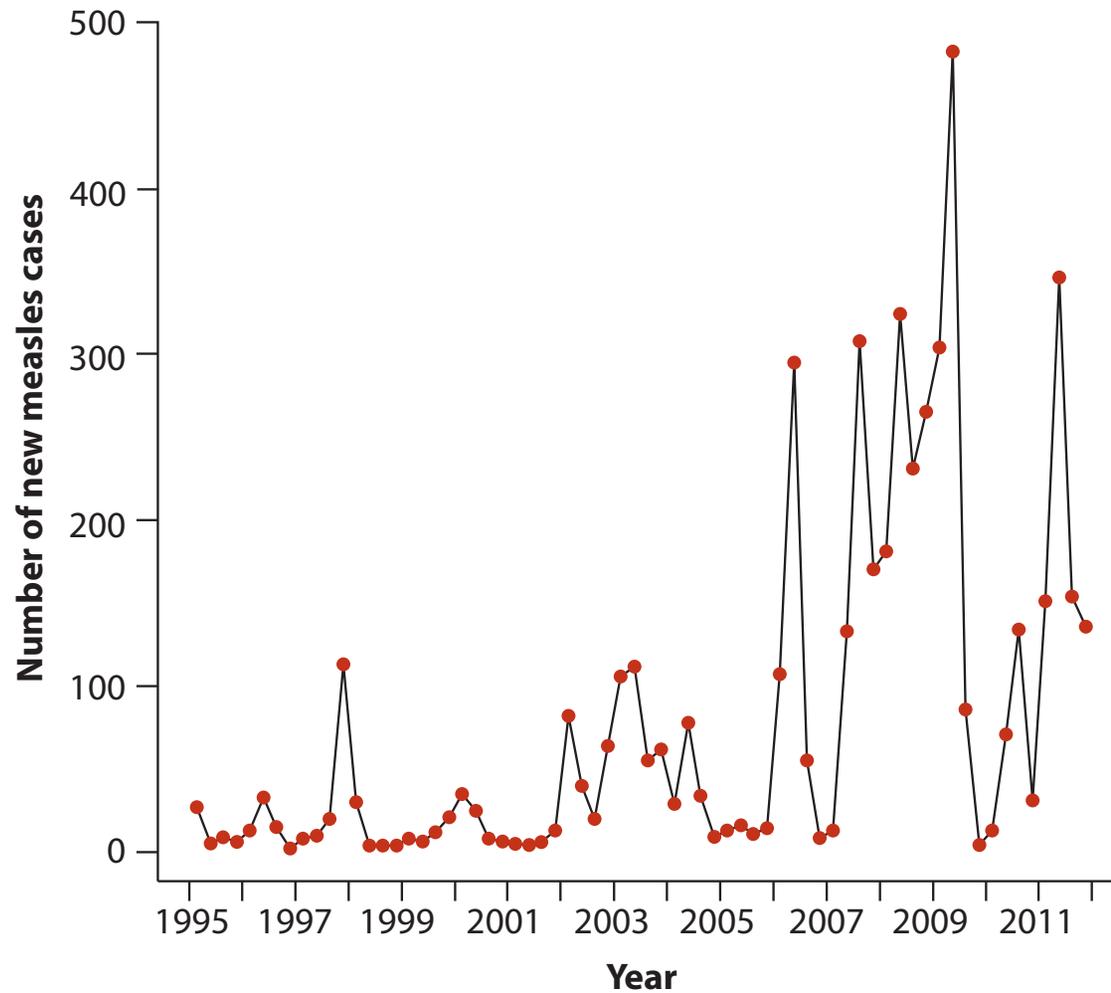


I. virginica

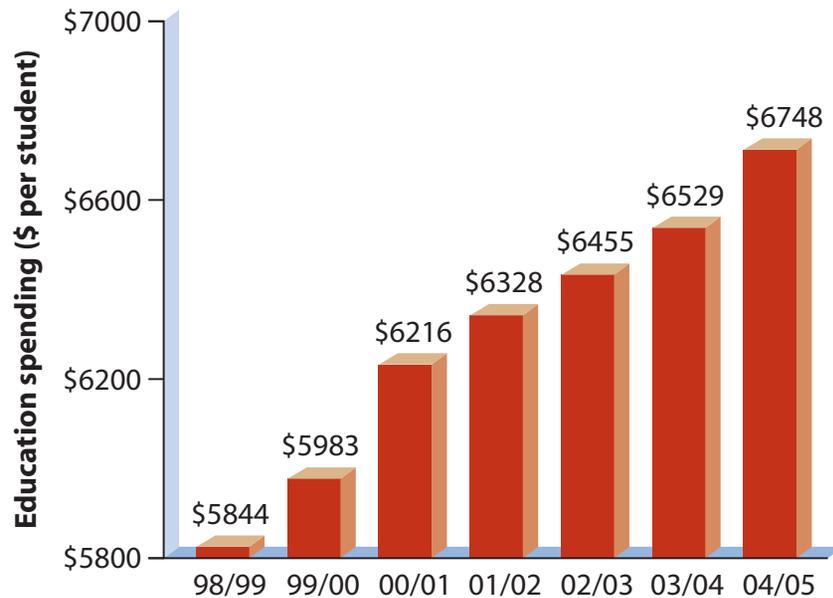


Sepal between petals

LINE GRAPH: uses dots connected by line segments to display trends in a measurement over time or other ordered states (e.g., size, etc).

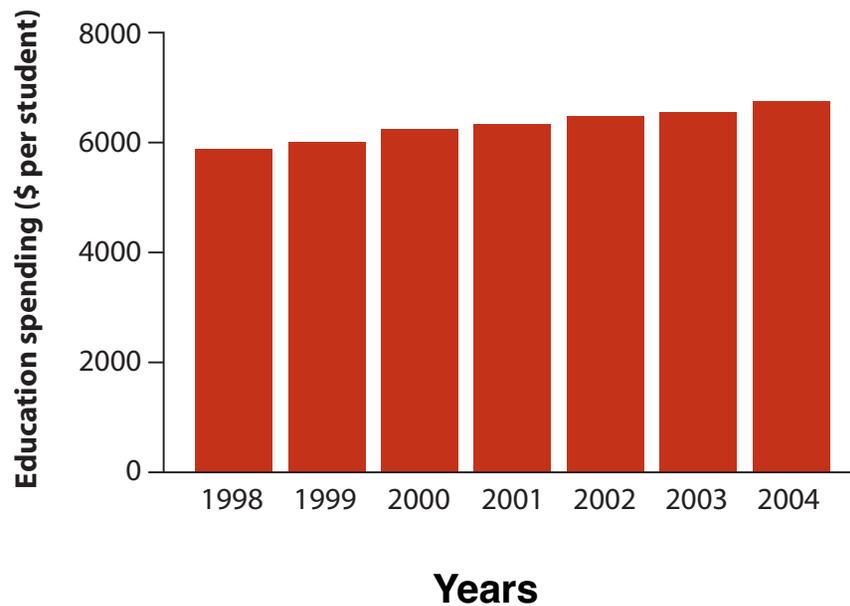
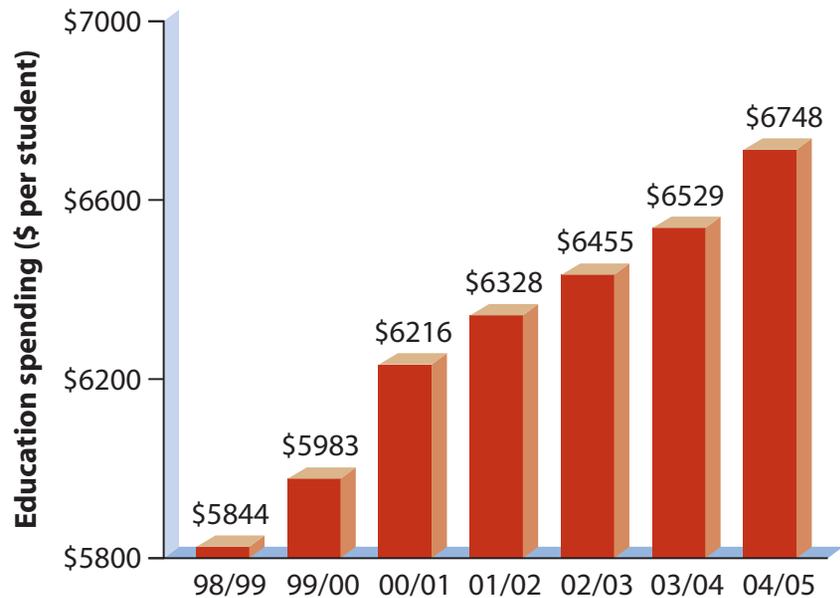


Graphs can be misleading! (both based on the same data): bad axis limits



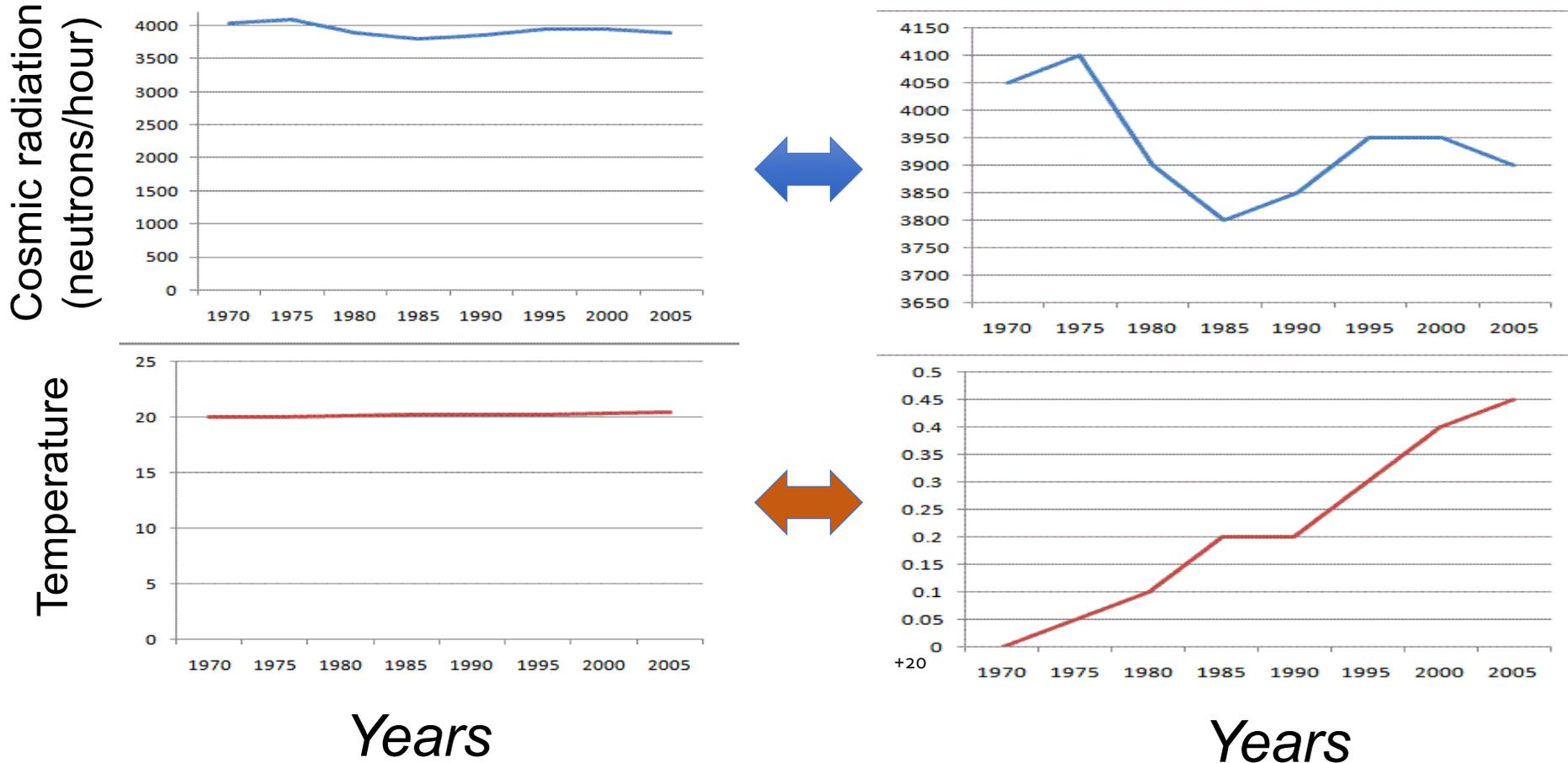
Upper graph. A bar graph, taken from a British Columbia government brochure, indicating education spending per student in different years. *Lower graph:* A revised presentation of the same data, in which the magnitude of the spending is proportional to the height and area of bars. This revision also removed the 3-D effects and the numbers above bars to make the pattern easier to see. The upper graph is modified from British Columbia Ministry of Education (2004).

Graphs can be misleading! (both based on the same data): bad axis limits



Upper graph. A bar graph, taken from a British Columbia government brochure, indicating education spending per student in different years. *Lower graph:* A revised presentation of the same data, in which the magnitude of the spending is proportional to the height and area of bars. This revision also removed the 3-D effects and the numbers above bars to make the pattern easier to see. The upper graph is modified from British Columbia Ministry of Education (2004).

Report the “intended” interpretation!



Axis limits do not change the data themselves, but they can strongly influence how the data are interpreted.

Some “rules” of Data visualization

How to make a *bad* plot:

1. Hide the data.
2. Make patterns hard to see.
3. Display magnitudes dishonestly.
4. Draw graphics unclearly.

How to make a *good* plot:

1. Show the data.
2. Make patterns easy to see.
3. Display magnitudes honestly.
4. Draw graphics clearly.

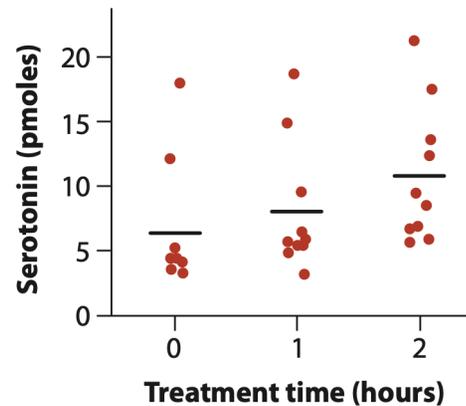
Mistake 1: Hide the data

How to hide data:

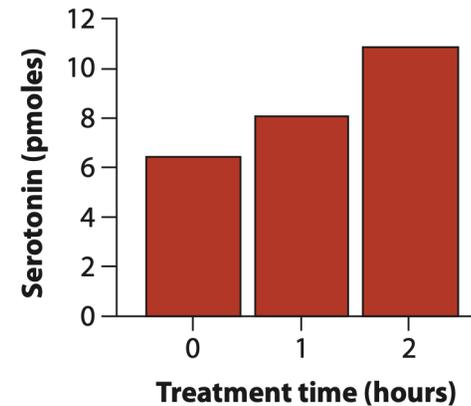
- Provide only statistical summaries (e.g., means).

How to reveal data:

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



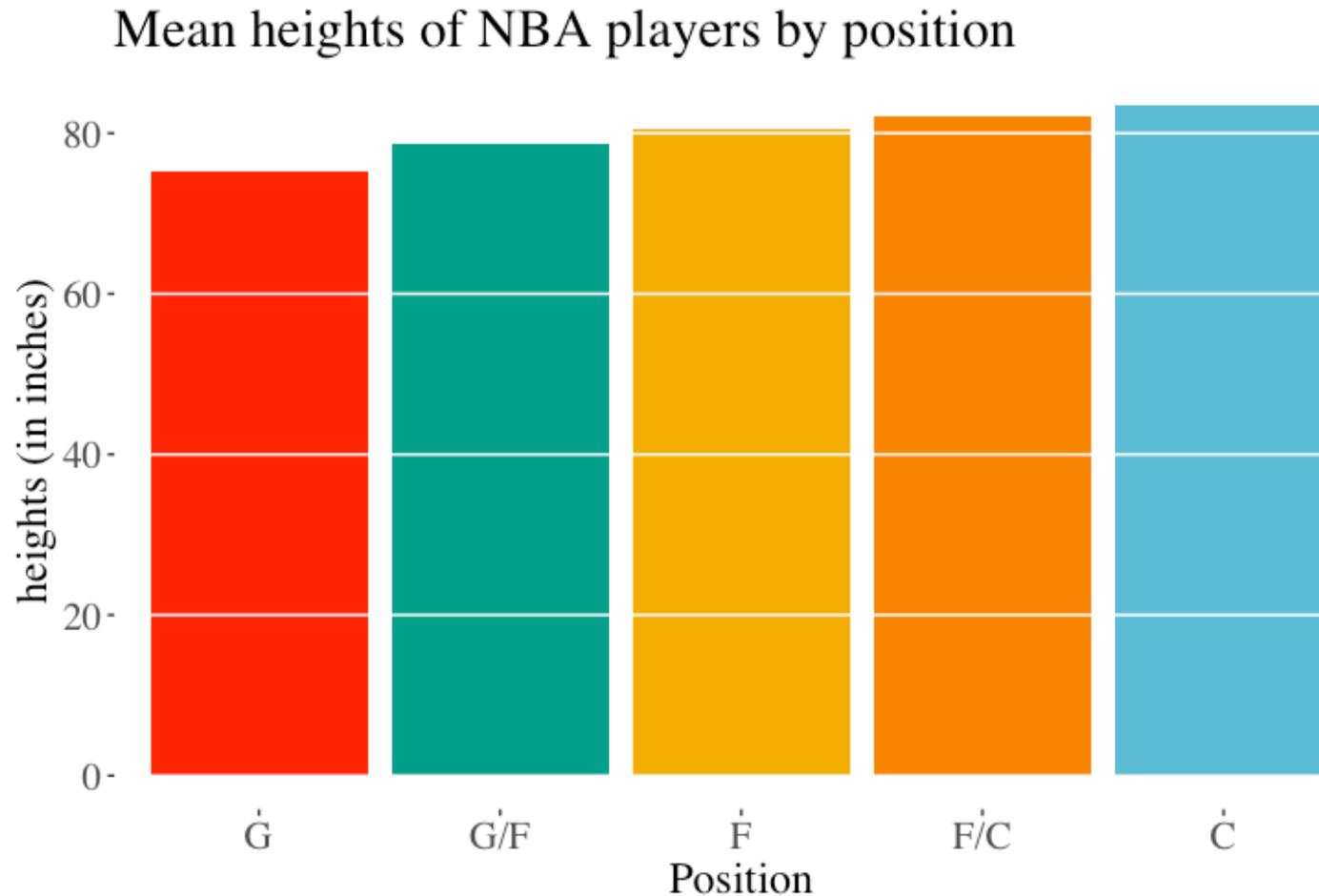
lines represent mean values
dots represent the data



bar heights represent
mean values

Not Showing Data, Just Summaries

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

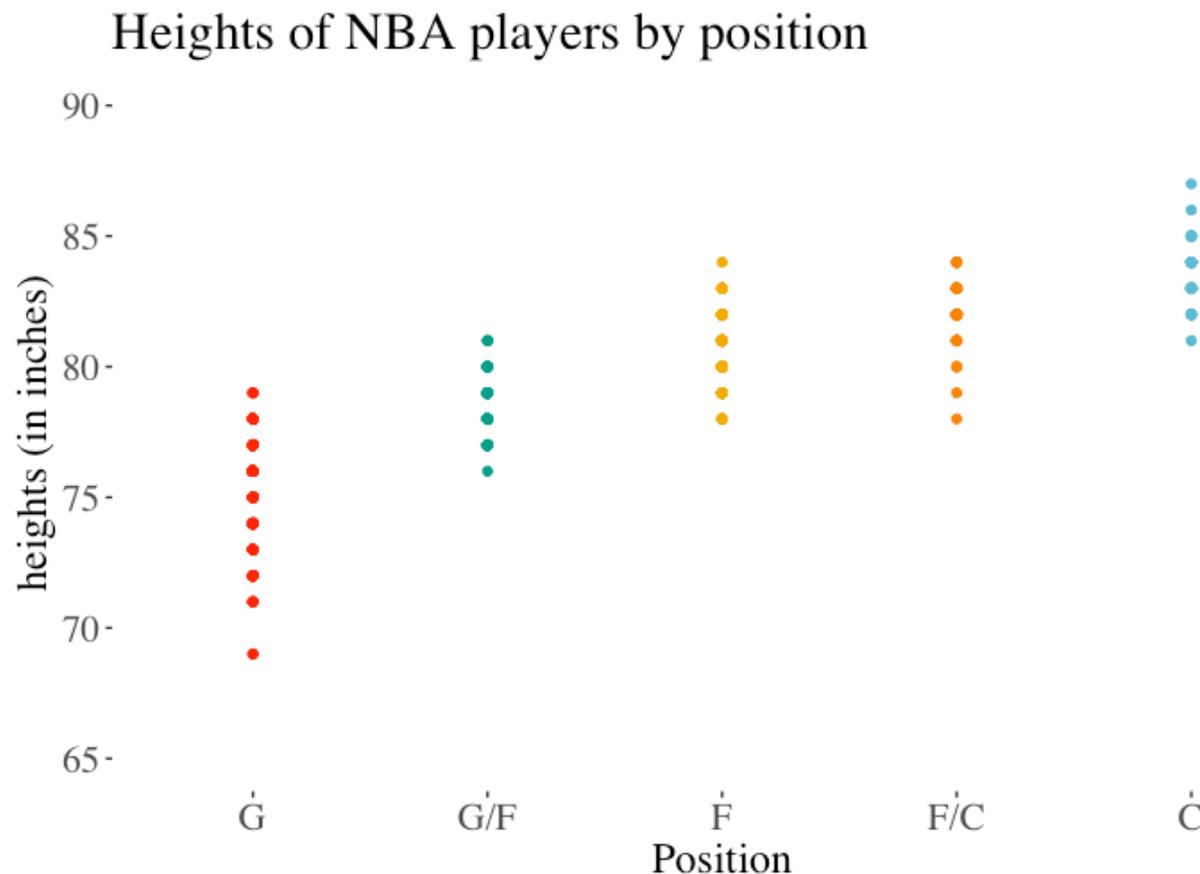


center, forward, guard

© 2020 W.H. Freeman and Company

Not Showing Data, Over-Plotting

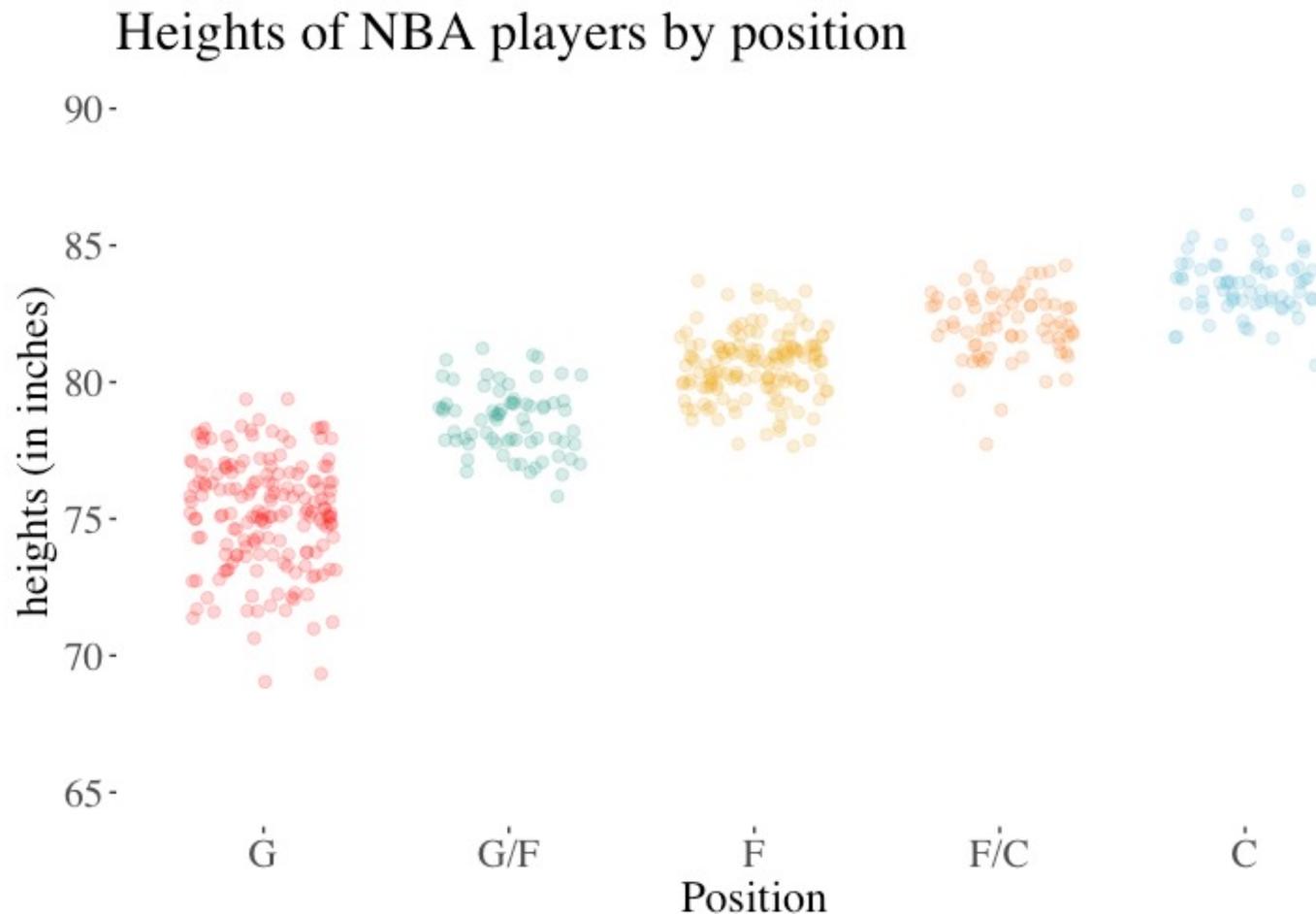
This plot hides the density of observations (number of players, i.e., many have very similar heights).



center, forward, guard

Showing Data, Jittering

This plot shows all the observations (all players and their heights).



center, forward, guard

© 2020 W.H. Freeman and Company

Mistake 2: Making Patterns Hard to See

How to hide patterns:

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

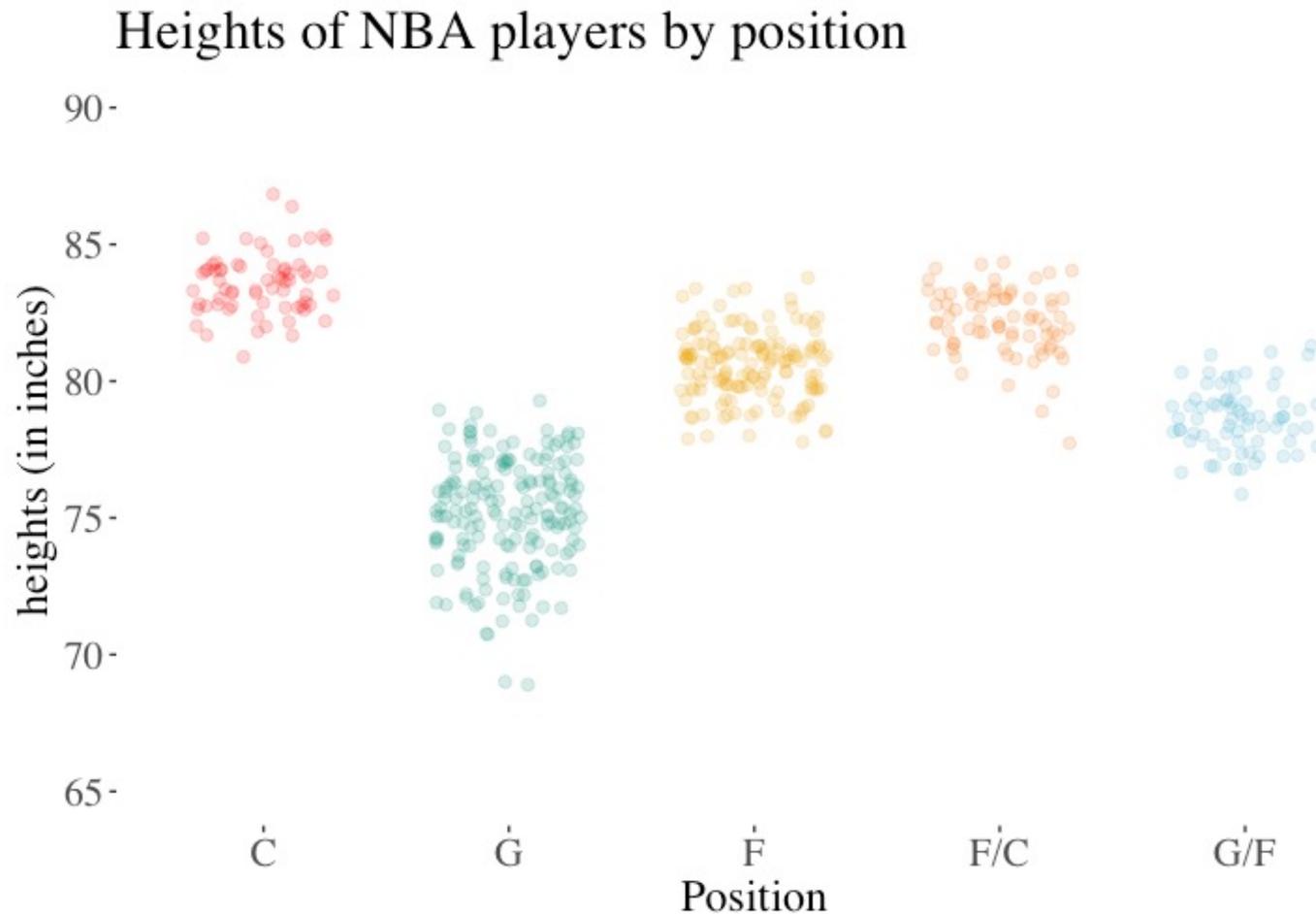
How to reveal patterns:

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

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

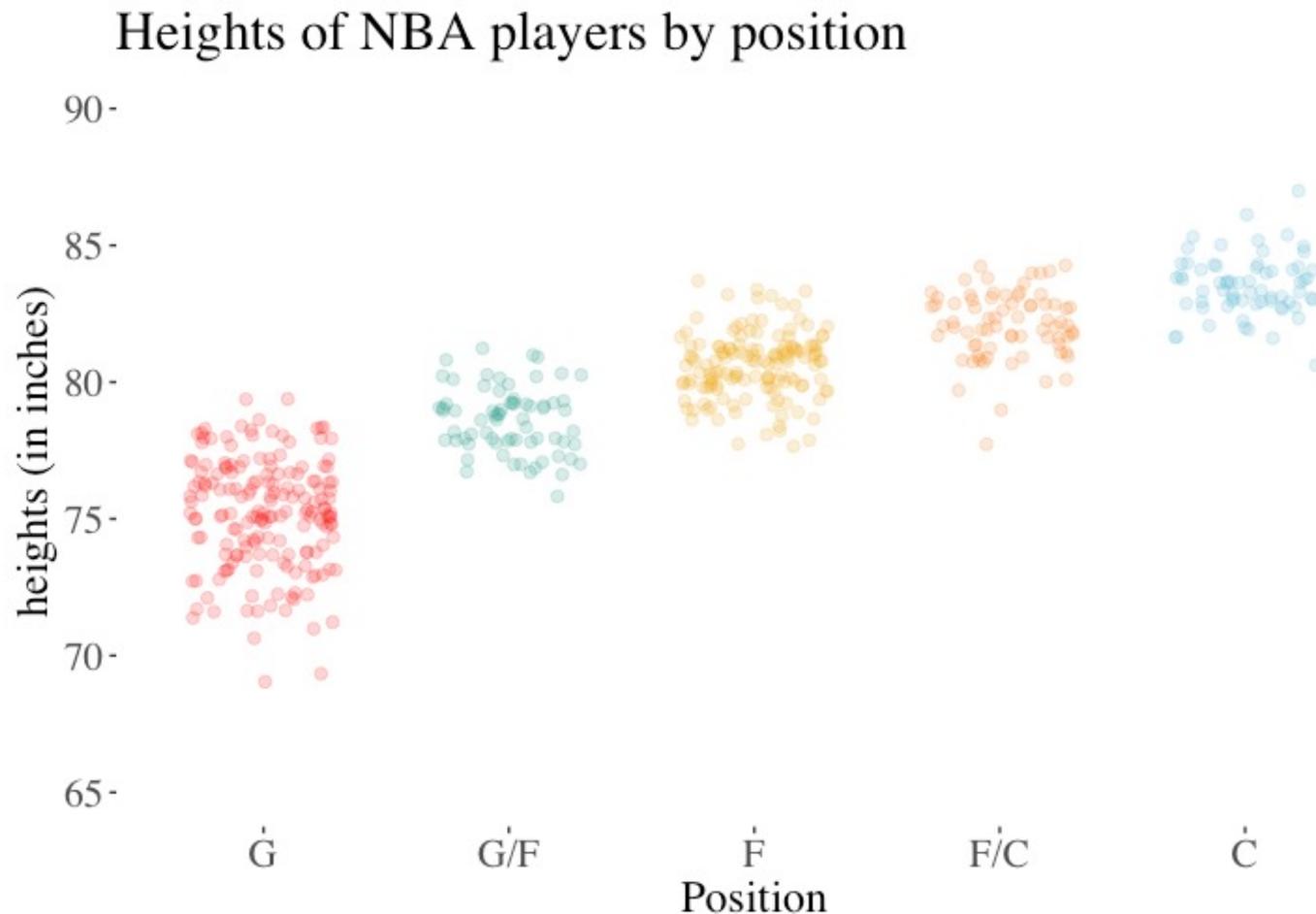
Nonsensical Order Hides Patterns

Nonsensical ordering of factors hides patterns.



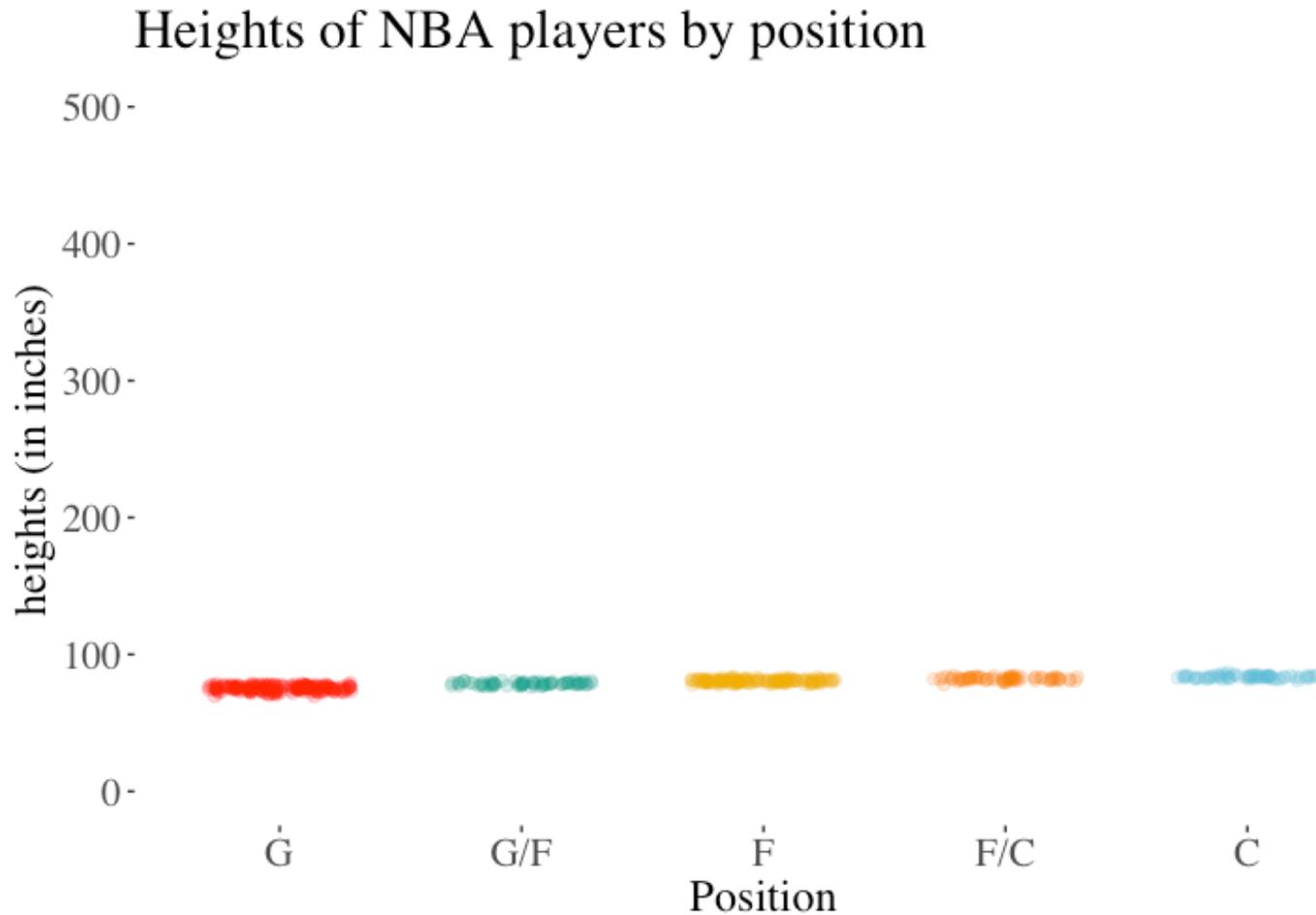
Nonsensical Order Hides Patterns

Sensical ordering of factors make patterns more evident (i.e., arrange by mean for nominal).



Bad Axis-Limits Hide Patterns

In this plot, the large scale (limits of the Y-axis) hides the pattern.



Graphs: The art of designing information

“A picture tells a thousand words”

- *Lake Blanche*

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

