

What is statistics?

Statistics is the science of learning from data.

Let's get some data: Survey: what is your favorite color.

red, orange, yellow, green, blue, purple, black, brown, grey, white other, e.g. pink, maroon, magenta,

Cases are the objects described by a set of data.

What are the cases here?

Students present in stat 202 today,

This was confusing because we collected data as a summary instead of the full data set

red 3	} instead of	{	bill → red
blue 4			sally → blue
white 2			pat → red
	erin → white		

A variable is a characteristic of a case.

What are the variables here (in the full data set)

name and favorite color

A label is a special variable used in some data sets to distinguish cases.

Each case must have a unique label (otherwise it is not called a label).

What labels are in full data set

name.

But two students can have the same name, especially first name) If we are collecting data we might want to use a different label

Student number, SSN, phone number, email address, made up number (e.g. 1 to 30)

Types of Variables

categorical variables - a variable that places a case into one of several groups or categories

quantitative variable - a variable that take numerical values for which arithmetic operations such as adding and averaging make sense

Favorite colors? which is it? (categorical)
Weight in pounds of students in class

Age: quantitative
Height: quantitative

Result of following survey:

I like pink: Strongly agree
agree
disagree
Strongly disagree

Party Affiliation: Green, Democrat, Republican, Libertarian
Unaffiliated

Gender Identity: male, female, Both, Neither

Another distinction is made between two different types of categorical variables

Ordinal-categorical variables have a natural order

Nominal-categorical variables are categories in name only, lack a natural order

Strongly agree
agree
disagree
strongly disagree } have a natural order

Colors: pink, grey, black, red are nominal categorical, because they lack a natural order

You can assign numbers to categories
But if the numbers are arbitrary
adding and averaging won't make sense

Consider grades A, B, C, D, F, A⁺, A⁻, B⁺, B⁻
Ordinal categorical

Or should we consider grades as quantitative, after all we compute GPAs,

To the extent that we compute GPAs ~~then~~ grades are quantitative. But many people argue that averaging for GPAs doesn't make sense. Is a F and a D or a D and a C the same distance from each other as a A and a A?

They argue that grades should be considered as ordinal categorical and care should be taken in computing GPAs.

~~When~~ Active Learning I

Problems 1.14 and 1.16
on Homework #1

key concept

Pg 6

The distribution of a variable tells us what values it takes and how often it takes these values,

Exploratory data analysis - when we study a data set, we

(1) ~~study~~ Study each variable by itself then move on to

(2) study the relationships among variables

For both of these steps we

- (a) Start with a graph or graphs
- (b) Then add numerical summaries.

The graphs chosen depend on whether the variable(s) are categorical or quantitative.

For single variables that are categorical there are two options

- * bar ~~graph~~ ^{plot}
 - * pie chart
- } both spins distributions (see def'n)

StatCrunch

* pie chart with summary
with data

↳ Favorite colors

It is important that the pie
add up to a whole.

In other words if there are too many
color lump the smaller ones into other

* bar ^{Plot} ~~graph~~ with summary
with data

No longer need to add to whole

* Loading files

* Value ascending etc

Homework: 1.22 - 1.27

Exploratory data analysis

Graphs for single quantitative variables

Stemplot } again shows
Histogram } distribution

Stemplots

Homework #2

Explain and do.