

Review

cases - objects described by data

variable - a characteristic of a case

label - variable(s) used to distinguish cases

categorical variable - places case into one of m several categories

quantitative variable - a variable that takes numerical values for which arithmetic operations make sense,

ordinal categorical - natural order

nominal categorical - in name only

Distribution of a variable tells us what values it takes and how often it takes them

Exploratory Data Analysis

I) Study each variable by itself

II) Study relationships among variables

For each: (a) start with a graph or graphs

(b) add numerical summaries

graphs for categorical: bar plot, pie chart } distribution
} show both

Pg 2

Graphs for
Quantitative variable

stemplot histogram

See homework 2; pass out extra copies
Project → if necessary

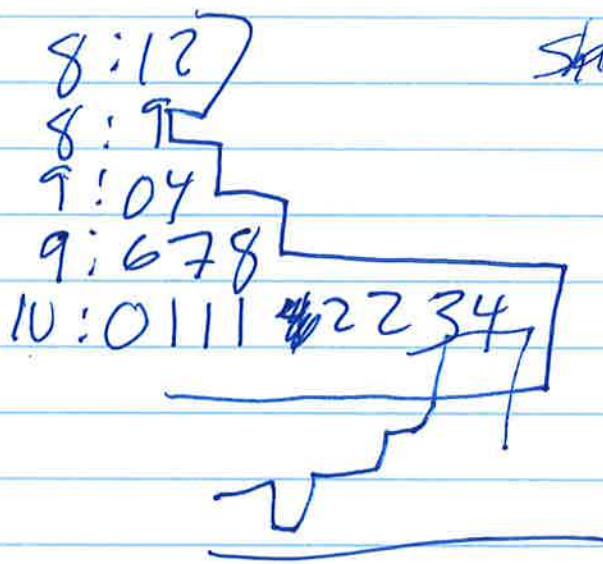
Do 1.7.

Pg 3

~~Pg 9~~

The other way of plotting the distribution of a single ~~categorical~~ quantitative variable is with a histogram

You can make a histogram from a stem plot as follows



slide in and rotate

bin with 5

pass out
homework
#3

Histograms are conventionally drawn vertically and not horizontally

Stemplots can be made with either by hand or software

Histograms should be made with software

Pg 4

Pg 10

Stemplots preferred when you have a
little bit of data

Histograms preferred (necessary when you have lots)

10K data point on stemplot?

Stemplots have bin widths that are
either 2, 5, 10, 20, 50, 100, etc

2 and 5 divide 10

Histograms can have any bin width

The counts of observations in bins are
called frequencies

80-85	2	}
85-90	1	
90-95	2	
95-100	2	
100-105	8	

A histogram is essentially a bar ~~graph~~ plot
of the freq
table EXCEPT that customarily we do not put
spaces between bars in a histogram (as we do for barplot)

Pg 1

~~Pg 11~~

Homework 3 CO₂

Load and plot histogram

No time yet to work on homework

#3

~~Pg 6~~

~~Pg 12~~

The height of a bar can be
(in StatCrunch)

frequency (# of datapoints in bin)

relative frequency $\left(\frac{\text{frequency}}{\text{total # of data points in set}} \right)$

density $\left(\frac{\text{relative frequency}}{\text{bin width}} \right)$

For these choices the shape of the histogram doesn't change only the vertical scale

Choice of bin width

No one right choice

too large (all data fall into one or a few bins)
"Sky Scraper"

too small (each datum falls into its own bin, each bin has only one data point) "pancake"

neither are desirable - use your best judgement

Note: changing bin width changes shape of histogram

Pg 7

~~Pg 13~~

Another difference between bar graphs and histograms

A bar plt need not have any measurement scale on the horizontal axis (if the variable is categorical) It simply identifies the item's being compared.

A histogram's horizontal scale indicates possible values of a quantitative variable

No space between bars is suggestive of the fact that all values are covered.

The horizontal scale of a histogram usually has chits.

~~Workshop Call Method~~

~~Tomorrow we will look at call center free to play around with it~~

Homework #3

We have data on the length of all 31,492 calls made to the customer service center of a small bank in a month.

Load and look at call center data set,

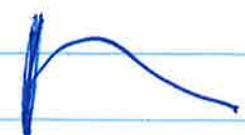
How do you analyze a dataset like this?

- Understand the background of the dataset (cases, variables, units of measure) (calls, length of call, seconds)
- Plot data
- Look for an overall pattern

bin width 10

where length $\leq 100 \leq 1000 \leq 5000 \leq 2000$

- Look for striking deviations from the pattern

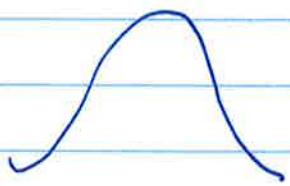
bizarr  what is called a long tail (lots of individual values much greater than most others)

explain 1st bin

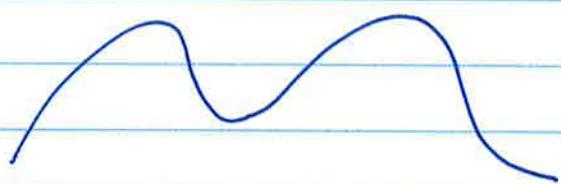
- Describe the overall pattern in terms of
 - shape
 - center
 - spread

Shape

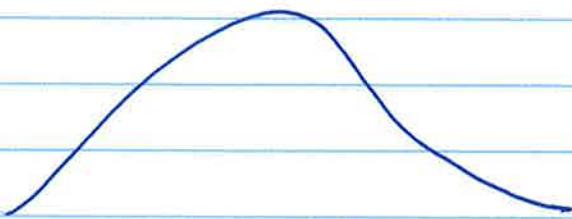
- the number of peaks or modes
- symmetric or skewed



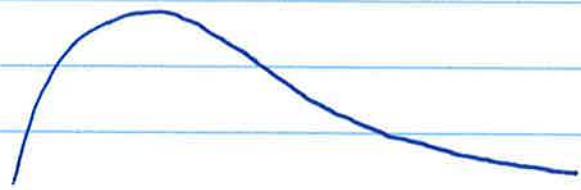
1 peak



two peaks



Symmetric
(about mid point)



Skewed to right

Call center two peaks (bimodal) skewed
Load IQ (one mode symmetric)

Skewed to left +
(analogous)



Tail - extreme values are said to be
within the tails of the distribution

Center - "midpoint" of distribution

Spread - How "wide" or "narrow" the distribution is

In a minute we'll make these notions precise

Recall after we make graphs and describe in terms of Shape, center, spread we add numerical summaries

Shape is more qualitative

for center and spread we ~~can't~~ can derive precise numerical summaries

Start with Center: there are two commonly used measures of center

Mean and Median

Mean - to find the mean of a set of observations, add their values and divide by number of observations

Notation for mean $\rightarrow \bar{X} = \frac{x_1 + x_2 + \dots + x_n}{n}$

Σ -notation $\rightarrow \bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$ limits of sum are 1 to n unless specified

Pg 11

Median - The median is the midpoint of the distribution. Half the obs are less than the median, the other half are larger.

To find median

1. sort the obs from smallest to largest
2. IF n (number of obs) is odd
median is center

7 obs total $\frac{\text{iii}}{3}$ \uparrow $\frac{\text{iii}}{3}$ not always equally spaced
median

If n is even median is halfway between center two

6 obs total $\underbrace{\text{o} \dots \text{p}}_{\text{3 obs}}$ $\underbrace{\text{...}}_{\text{3 obs}}$
median

Show mean and median on StatCrunch for IQ and call center.

PGR

Mean is sensitive to outliers
Median is resistant to outliers

Let's say there are 999 datapoints
and all 999 are equal to 10

What's the mean? 10

What's the median? 10

Now suppose we add an outlier - a 1000th
data point equal to 1 trillion 10^{12}

What's the mean

$$\frac{10 + 10 + 10 + \dots + 10^{12}}{1000} = \frac{10^9 + 9990}{1000}$$

$$= 10^9 + 9.99 \\ \approx 1 \text{ billion}$$

What's the median

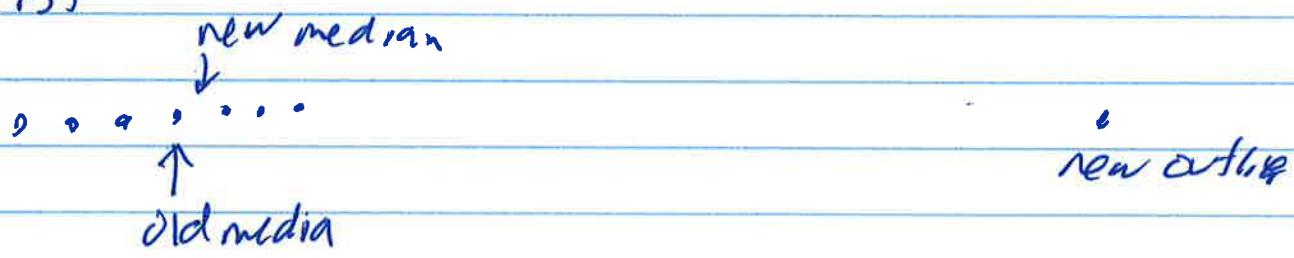
Middle points both 10 so 10.

Doesn't change at all

PG 13

One enormous outlier changes the mean tremendously but one enormous outlier usually does not change the median much. In this case not at all.

A single outlier can change the median by at most one data point in the sorted list



Median is resistant to outliers (or just resistant)

Mean is sensitive to ~~outliers~~ outliers.

Important distinction

If outliers are present it is ~~usually~~ often best to use resistant measures

Median and mean are measures for center of distribution. This doesn't tell whole story

- median family income doesn't tell you about extremes of wealth and poverty
- A drug may have correct mean concentration of an active ingredient but if some batches are way too high and others are way too low, that's a problem.

Need a measure of spread of a distribution

Spread (Resistant to outliers)

Upper Quartile

Lower Quartile

Min

Max

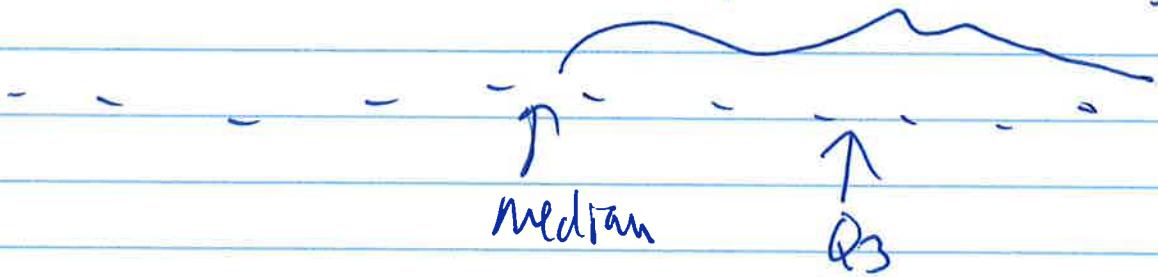
Spread (Sensitive to outliers) ← more commonly used
Standard deviation

Upper Quartile (aka Third Quartile) (denoted Q_3)

is the median of the upper half
of the observations (i.e. the median of
the observations above the overall median)

Also called the 75th percentile

median of these is Q_3



Lower Quartile (First Quartile)

Median of lower half of obs

Also called the 25th percentile

The p^{th} percentile - the value such that
 p percent of observations fall at or below it

Median is the 50th quartile percentile
Or 2nd quartile Q_2 (but this notation
is not usually used M is used.)

(Pg 16)

Alert for percentiles: There might not be an obs such that exactly p percent of obs fall at or below it.

Example 11 obs find 27th percentile

Book says: take nearest obs

Software: different programs might give slightly different values.

For 25th, 50th, 75th percentiles there is an exact def'n based on the def'n of median (middle obs or medpnt of middle two)

But even in this case software can give different answers

OK to report whatever software gives
DF publishing, say what software used.

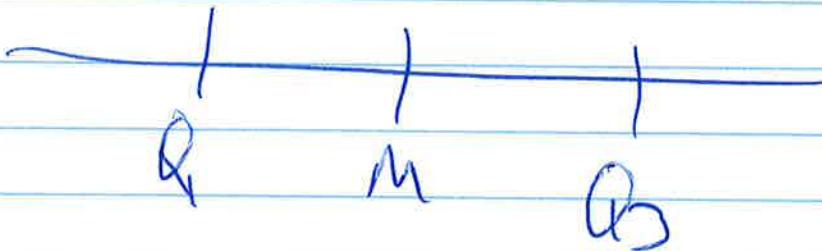
Quartile and Percentiles are resistant to outliers

Pg 17

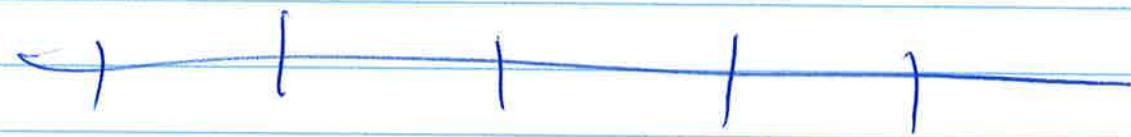
For describing center and spread of a distribution there is something called the 5-number summary

of a set of observations

Three of these numbers are



The other two are

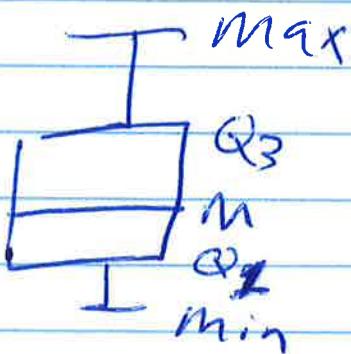


(Smallest
obs)

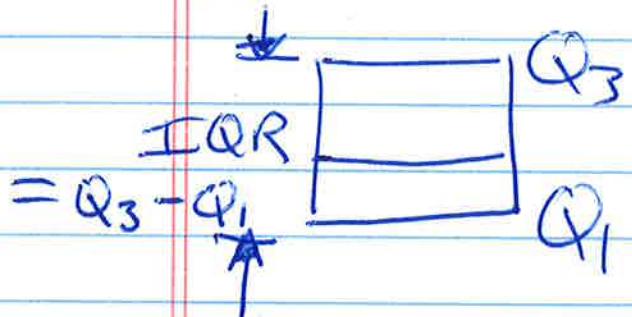
(largest
obs)

(Pg 68)

You can display the 5-number summary with a box plot



The distance between Q_1 and Q_3 is called the inner quartile range (IQR)



$$IQR = Q_3 - Q_1$$

The IQR is useful for flagging

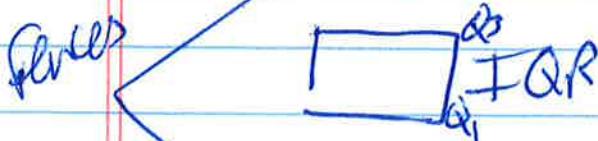
Potential outliers

1.5 IQR Rule for Outliers

Call an obs a suspected outlier if it falls more than $1.5 \times IQR$ above Q_3 or below Q_1

$$T \approx Q_3 + 1.5 IQR$$

points that fall beyond fences are flagged as potential outliers



$$Q_1 - 1.5 IQR$$

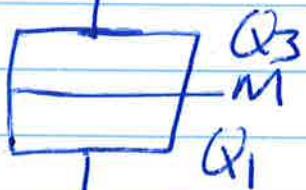
A modified boxplot also draws →
outliers

- * Don't include suspected outliers in min/max

- * Draw asterisks at suspected outliers

- * outlier

max excluding outliers



min excluding outliers

*] outliers

IQ
Show call center

Pass out and do homework #4

Talk about log transform