

AMS 7

April 3, 2017

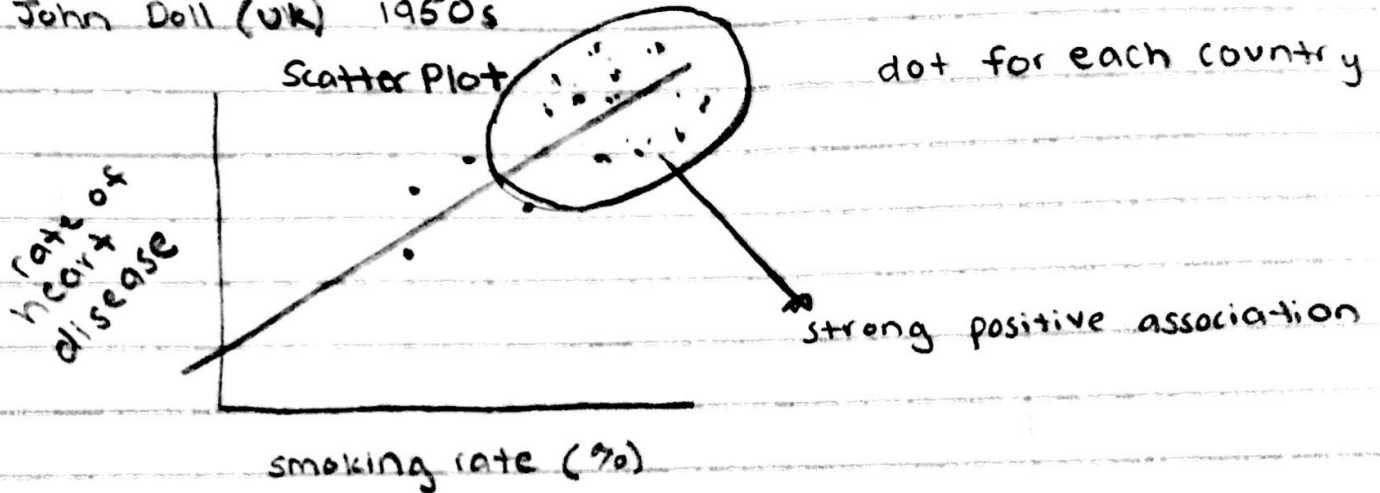
This time - intro

Next time - descriptive methods

Webcasts: webcast.ucsc.edu

↳ lectures will be found here

Dr. John Doll (UK) 1950s



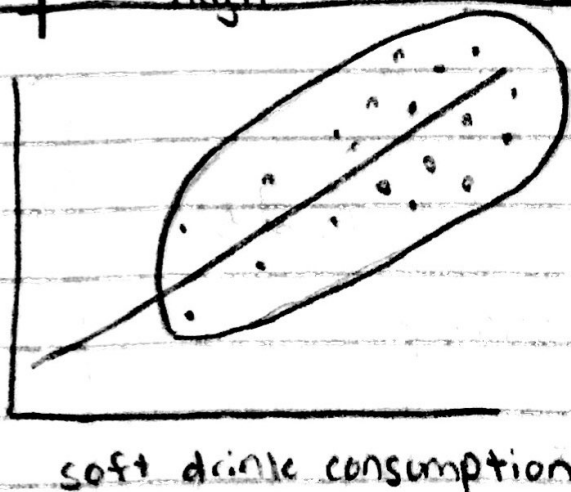
April 5, 2017

This Time - descriptive methods

Next Time - descriptive Methods

<u>Season</u>	<u>soft drink consumption</u>	<u>polio</u>
Fall	medium	medium
Winter	low	low
Spring	medium	medium
Summer	high	high

Polio incidences

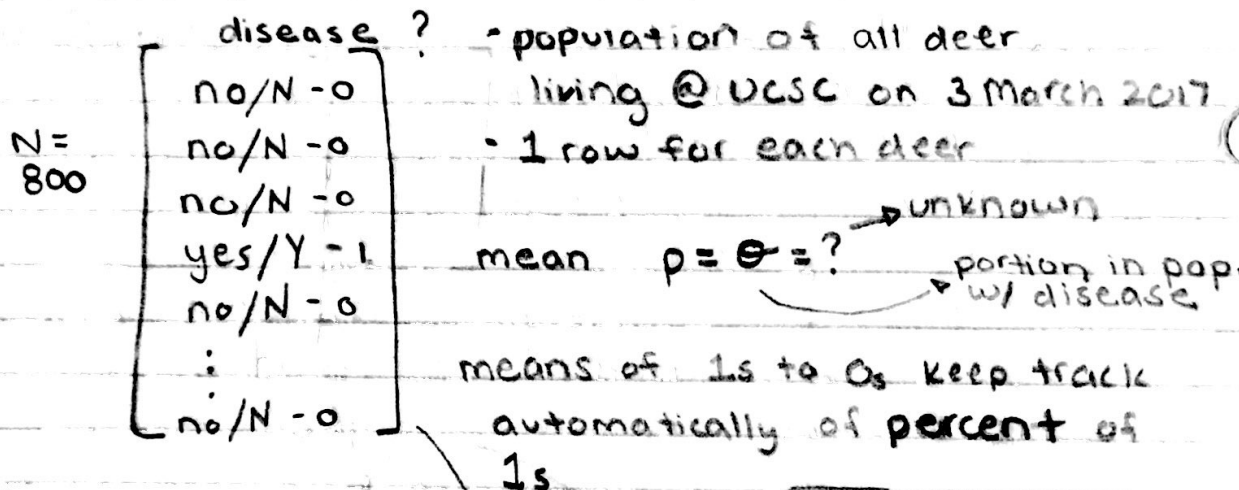
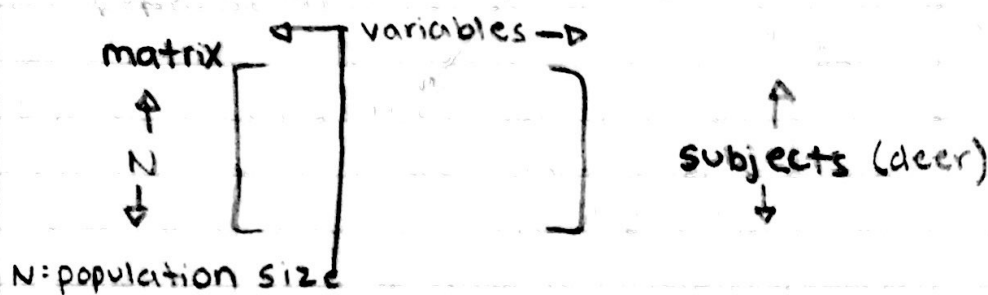


Lecture 1:

Stats - the study of uncertainty; how to measure it, and what to do about it

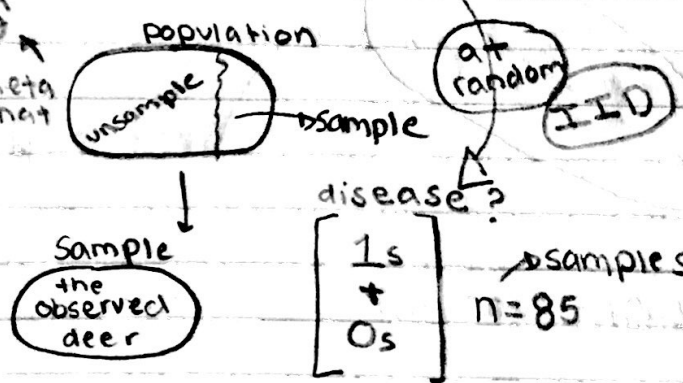
Uncertainty - state of incomplete or imperfect info about something of interest to you

Percentage θ (theta = p) \underline{P} = population



mean \bar{y} \hat{p} $\hat{\theta}$
 y bar \hat{p} hat θ hat

$n = 85$
 disease = 3
 $\frac{3}{85} = 3.5\%$ has disease



goal of sampling: representativeness: want sample and unsample to be similar in all relevant ways (IID)

- at random with replacement (Independent identically distributed) (SRS)
- at random w/out replacement (simple random sampling)

SRS - is more informative than IID, but
IID has easier math

if n is a lot smaller than N , SRS and IID are
approximately equal to each other

SRS \doteq IID

April 7, 2017

* Sections start Monday, April 10th

webcast

→ zero

username: ams-7-01

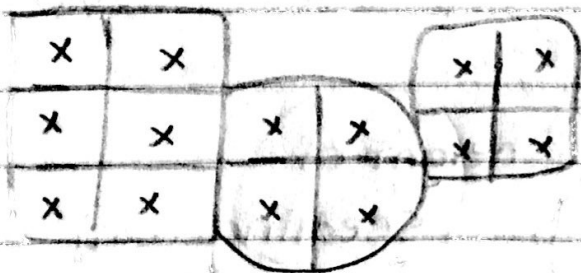
password: ready4slow

Lecture 2

- variables (things that can be measured on pop. subjects)

• variable can be either [yes, no]

- parameter (numerical summary of a population)



haphazard = random

Variables

possible values

- eye color in humans → brown, blue, black, green

- hair color in hamsters → brown, black, dichotomous

qualitative:
no unique place for its value on a number line

- success at maze running → very slow, slow, moderate, fast, very fast

ordered categorical or ordinal

qualitative:
order to category

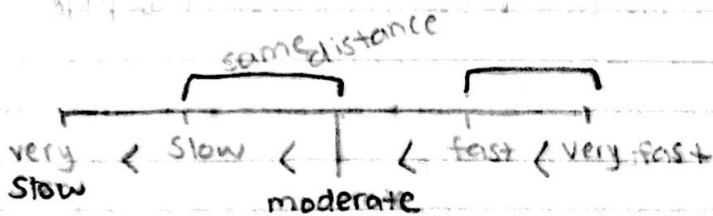
- seedling height (cm) → 7.42 cm, 7.017 cm

quantitative:

values have unique places on the number line

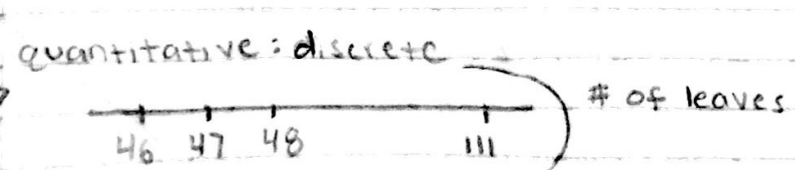
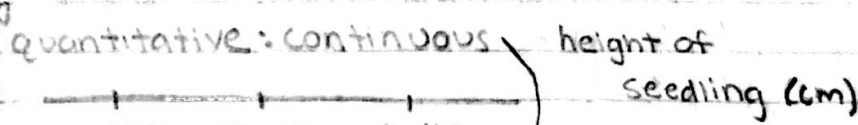
- # leaves (adult) → 46, 111

Plot size



ratio scale variable

has a true zero



constant size interval + ratio variable

Three variables measured: qualitative

Animal ID	Age	Hair Color
45	2	0
333	0	2
.	.	.
.	.	.
.	.	.

can it take mean of qualitative variables

0 = brown
2 = red

1 column for each variable

1 row for each animal

mean no point, not meaningful meaningful not meaningful

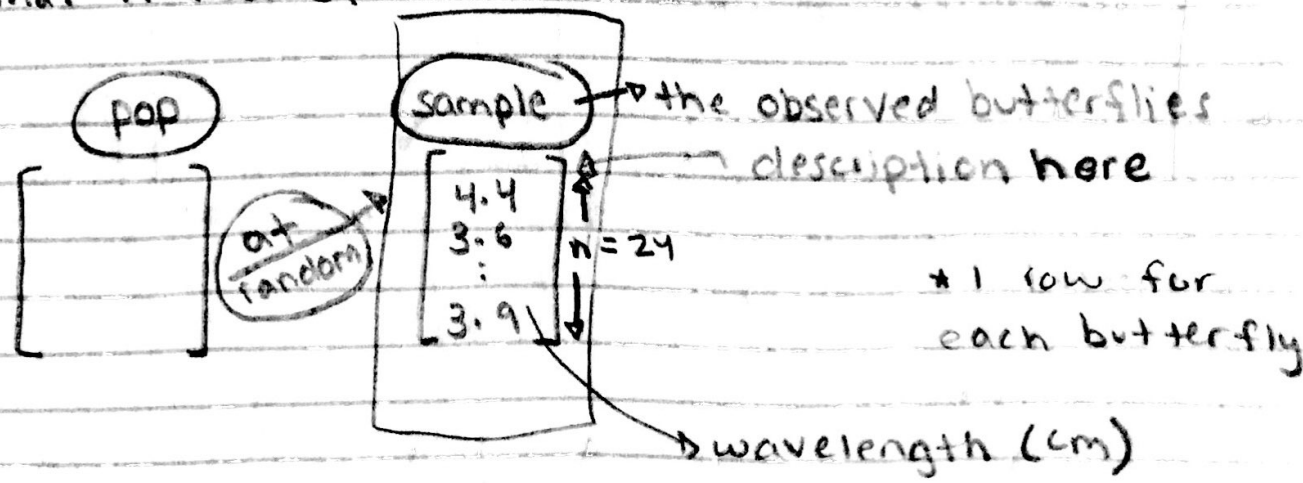
April 10, 2017

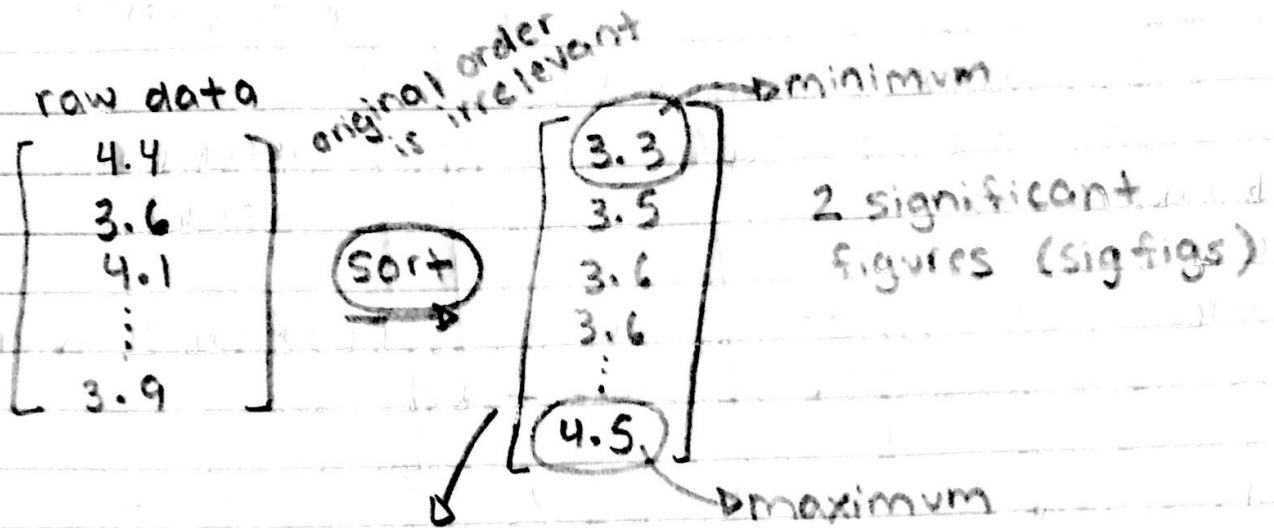
This time: histograms, bar graphs
Next time: measures of center spread

*discussion start this week, go! (quiz)

- I. intro (week 1)
- II. descriptive methods (graphical numerical)

you have a data set, how to gain insight to what it means?



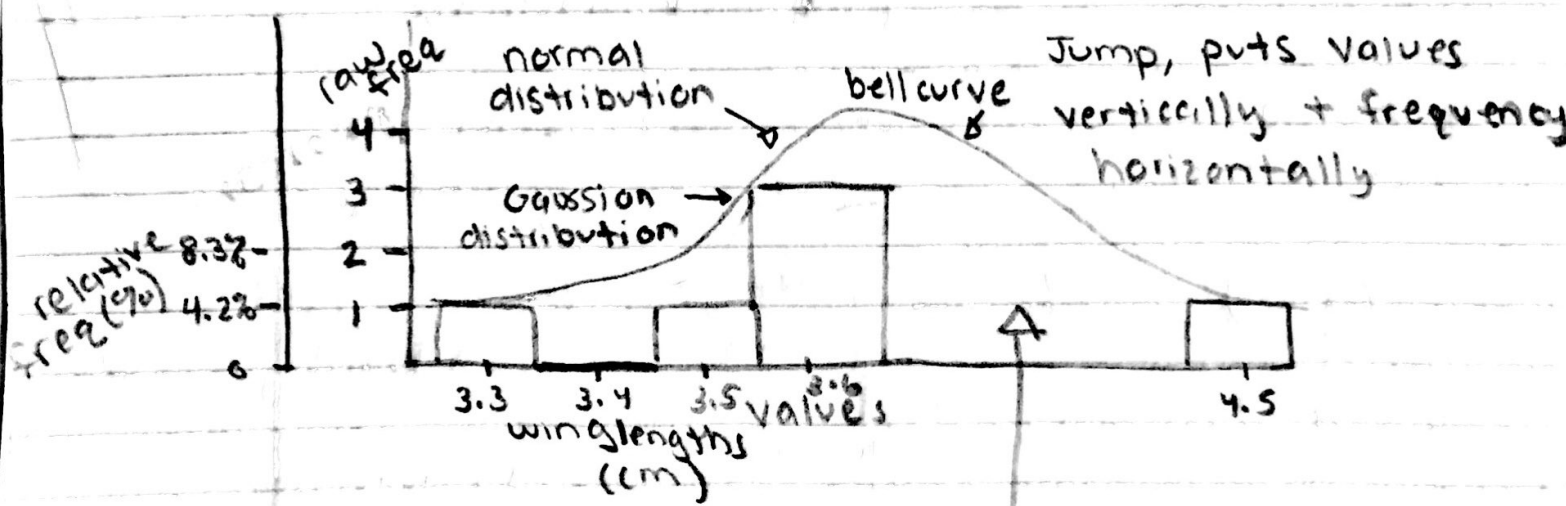


values	count = raw frequency
3.3	1
3.4	0
3.5	1
3.6	3
⋮	⋮
4.5	1

raw frequency distribution

n=24

variables are often named with values such as x, y, and z
 ex. y = wavelengths



* Carl Friedrich Gauss (1777-1855)

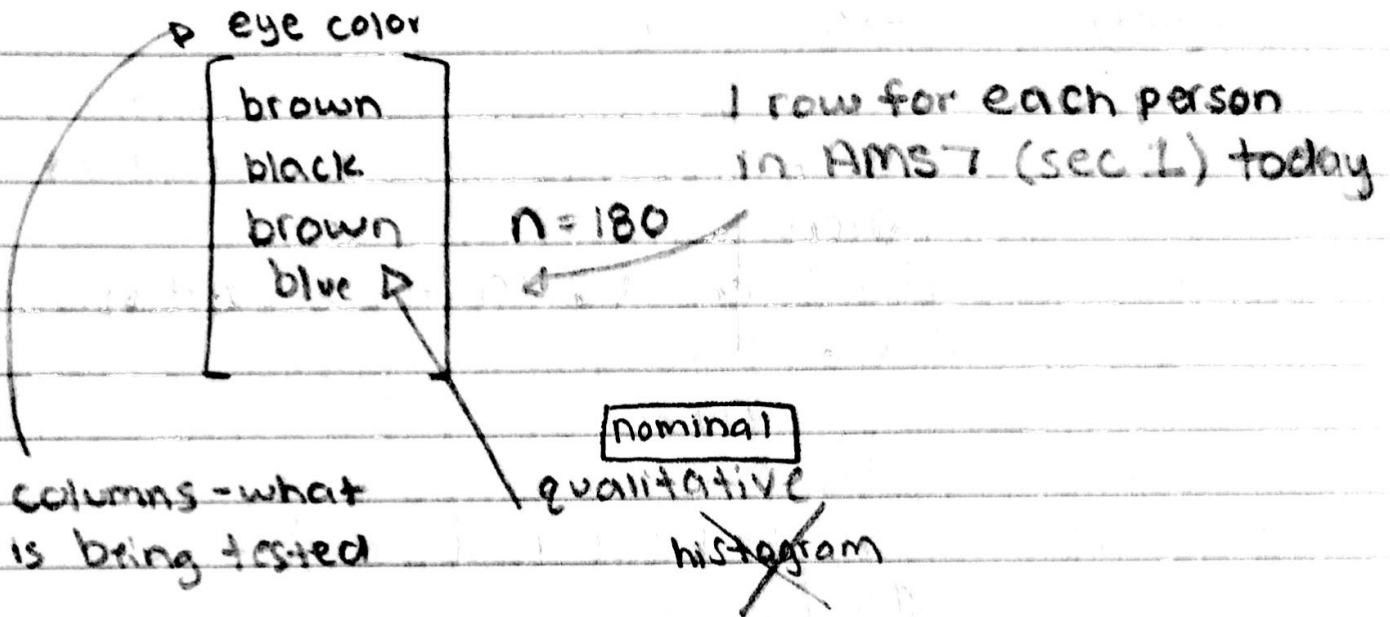
raw frequency histogram (a kind of bar graph)

wavelength ← quantitative

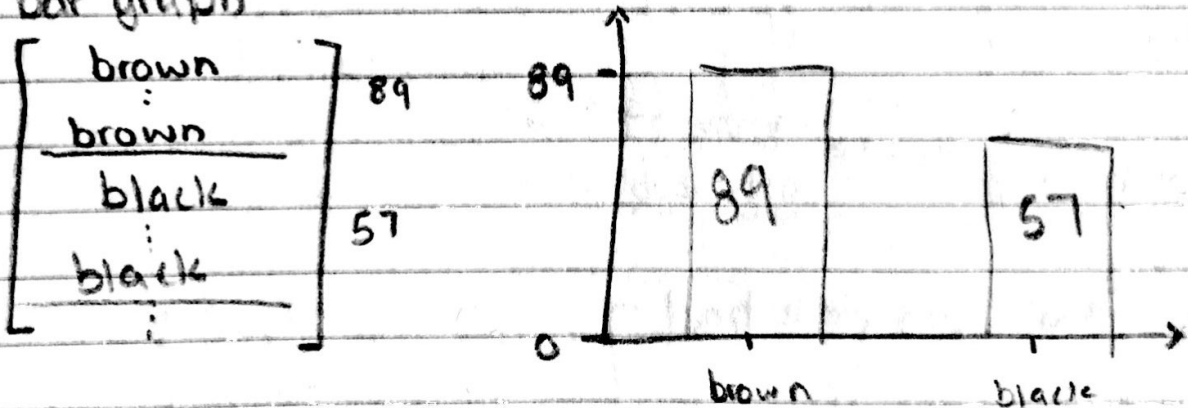
* conceptually continuous, but made discrete by rounding in measurement process

histograms: can only be made for quantitative variables (either discrete or continuous)

* visualize the raw data set



bar graph



bar graphs = qualitative

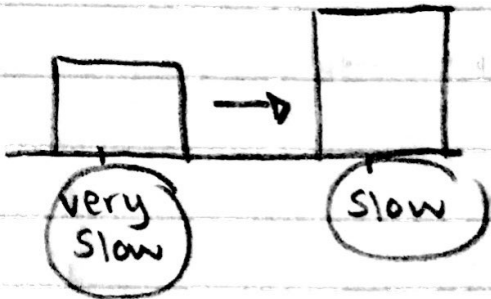
histograms = quantitative

value	raw freq.	relative freq (%)
3.3	1	$\frac{1}{24} \times 100\% = 4.2\%$
3.4	0	0%
3.5	1	$\frac{1}{24} \times 100\% = 4.2\%$
3.6	2	$\frac{2}{24} \times 100\% = 8.3\%$

* 3 diff. vertical scales for histograms

① raw frequency vs ② relative freq vs ③ density

maze running



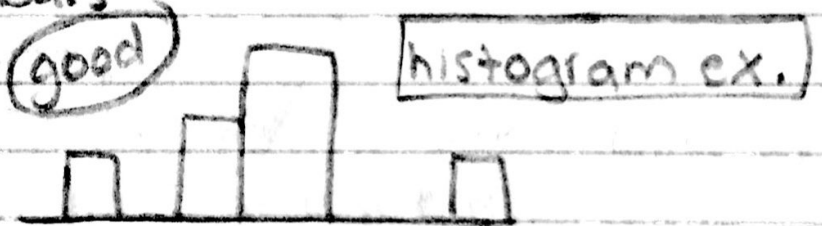
qualitative, ordinal

Question: Is there just one unique histogram for any given data set?

Answer: NO, you have to make a good choice of # of bars

butterfly: $n=24$

bar 13



bar

bad

another histogram

all sense of distribution shape is lost

3.3

4.5

too few bars

too few bars \rightarrow bad \rightarrow all sense of distribution shape is lost

just right \rightarrow good

too many bars \rightarrow can't see the forest for the tree (too much detail)

* to cut # of bars in half, combine adjacent bars

location

vine
building eave
vine
⋮

sparrows

$n=211$ 1 row for each nest (nesting site)

qualitative, nominal, non dich.

bar graph = yes
histogram = no

pigment type

0
2
5
0
⋮

sunfish

$n=$ 1 row for each sunfish

qualitative = even when we code 0, ... 4

bar graph = yes
histogram = no

April 12, 2017

This time: numerical measures of center spreads
Next time: the normal curve

ex. litter size in foxes

of pups

4
7
3
4

1 row for each litter
quantitative, discrete, ratio
- pups between possible values represent impossibility

histogram? yes
mean? yes

ex. aphids on clover

of aphids
14
36
⋮
⋮

1 row for each clover plant
n=424

histogram? yes
mean? yes

- quantitative, discrete, ratio

ex.

phosphorous concentration (y)

8.42
9.08
⋮
⋮

n=130

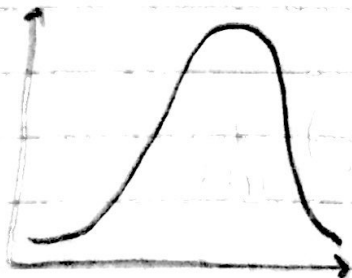
1 row for each leaf

- quantitative, conceptually continuous, ratio

histogram? yes
mean? yes

$$\frac{y_{\max} - y_{\min}}{\# \text{ bars}} = \text{bar width}$$

"law free"



value

relative freq. (%)

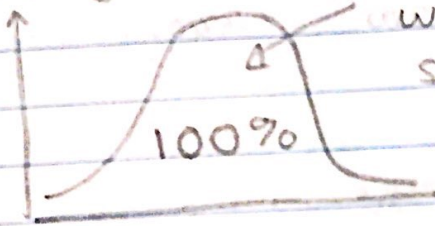
density scale

* Unless otherwise noted, all histograms in this class will be drawn on the density scale

density
Scale

$$\text{relative freq.} = \text{area under histogram (curve)}$$

density



with histogram on density scale, total area under curve = 1 or = 100%

histogram shapes

unimodal



phosphorous

symmetric point of

symmetry

barrier, no neg. income



income Trump

long right hand tail

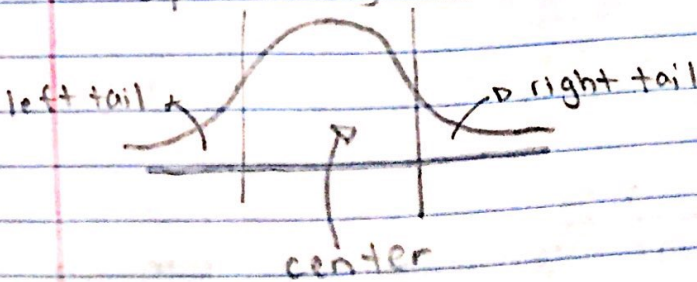
barrier, can't get higher than 100%



midterm scores

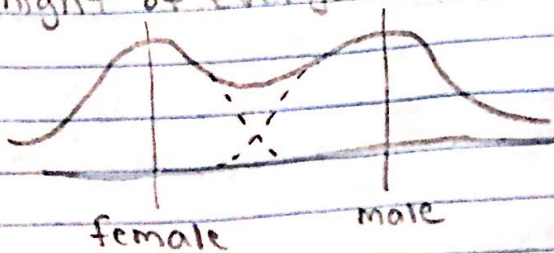
asymmetric = skewed

* long left-hand tail



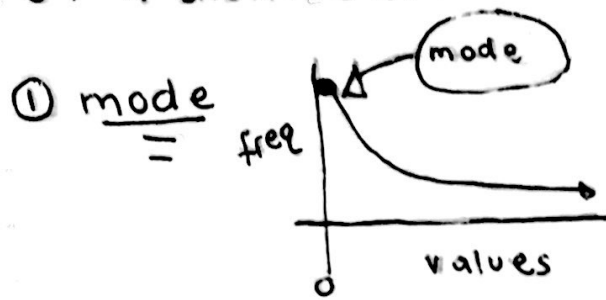
ex. height of everyone in class today

bimodal or multimodal

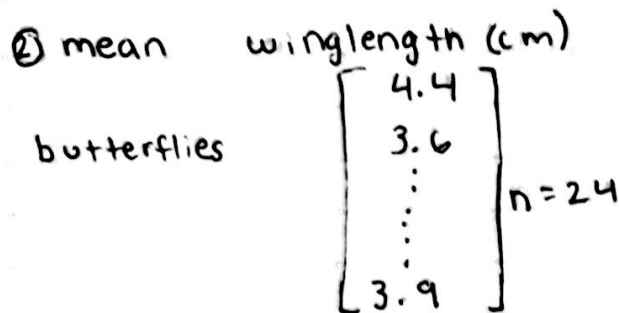


mode = 9.) a point of highest frequency

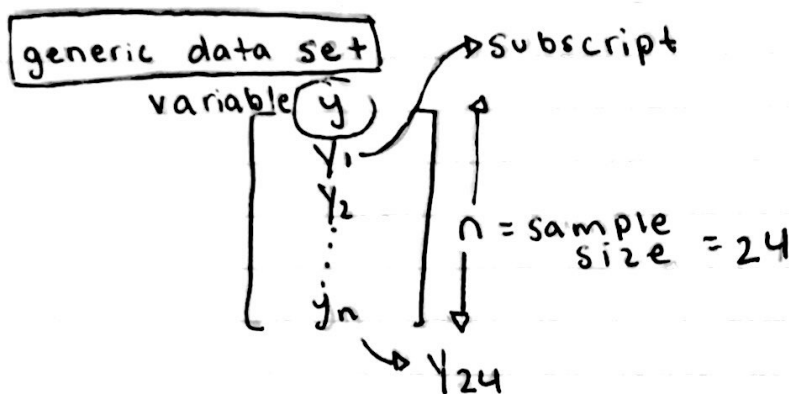
quantitative (numerical) measures the center of a distribution



mode is not necessarily central, but always typical in terms of frequency



mean = ?



mean \bar{y}
"y-bar"

$$\frac{y_1 + y_2 + \dots + y_n}{n}$$

$$= \frac{1}{n} (y_1 + y_2 + \dots + y_n)$$

capital sigma

summation notation

largest index

smallest index

index of summation (i, j, k)
(expand out theorem)

$$\sum_{i=1}^n y_i = y_1 + y_2 + \dots + y_n$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

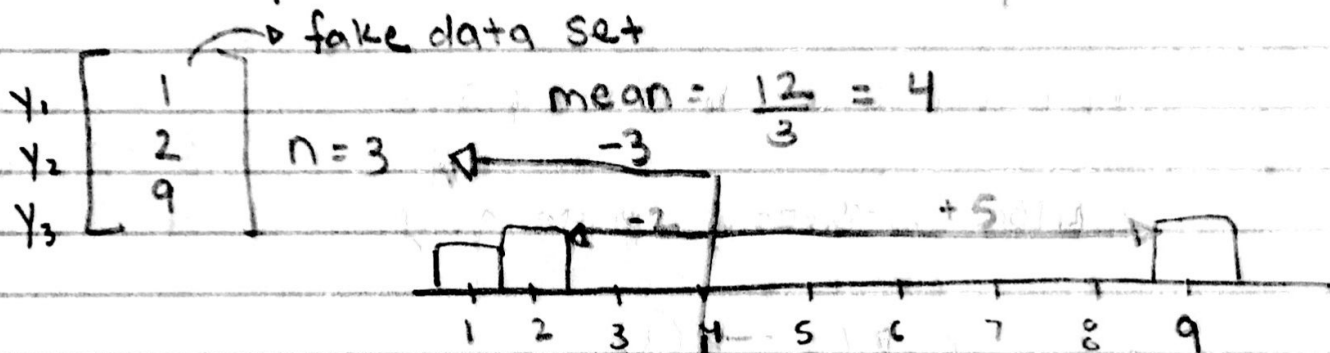
sample mean
of variable y

* in this data set

$$\bar{y} = \frac{95.0}{24}$$

$$= 3.96 \text{ cm}$$

graphical interpretation of the mean



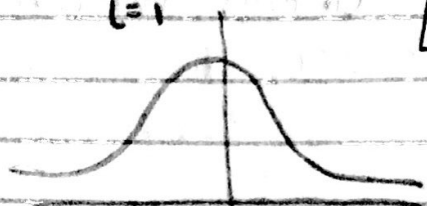
$\begin{bmatrix} 1 \\ 2 \\ 9 \end{bmatrix}$	$=$	$\begin{bmatrix} y_1 \\ y_2 \\ y_n \end{bmatrix}$	$\xrightarrow[\bar{y}=4]{\text{subtract}}$	$\begin{bmatrix} 1-4=-3 \\ 2-4=-2 \\ 9-4=+5 \end{bmatrix}$	$n=3$
---	-----	---	--	--	-------

theorem:

let n be any positive integer y_1, \dots, y_n be any real numbers, then

$$\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})$$

mean = balance point



symmetric +
unimodal

point of symmetry = mode = mean =

April 14, 2017

This time: standard derivative, the normal curve
Next time: controlled experiments + observational studies

Read: DD
ch 1-3 (A) , 1-6 (B) , LN pp. 1-84?

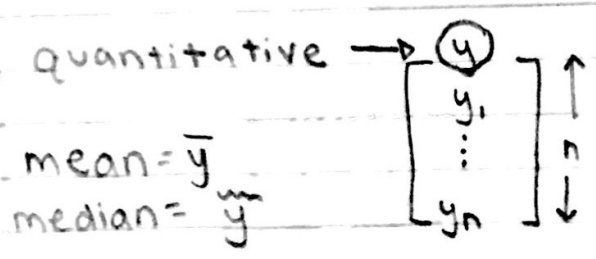
No more ecommons → HW is to be submitted to canvas.vcsc.edu

Official Note taker = Kaitlyn Abercrombie

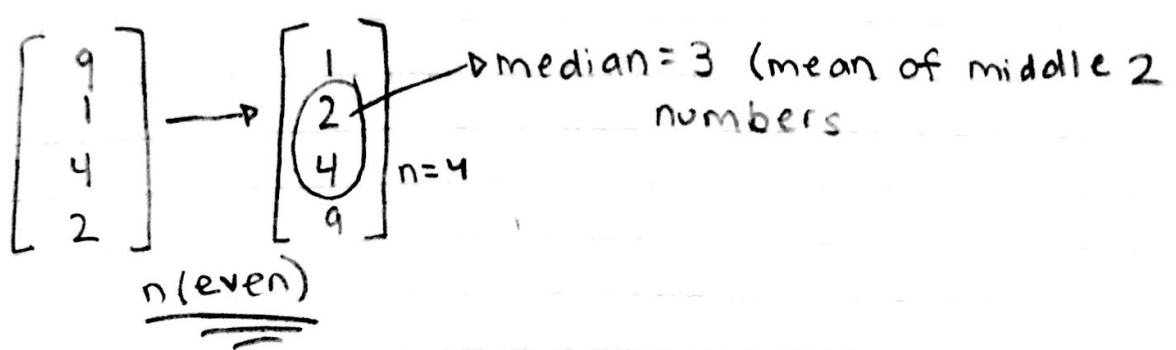
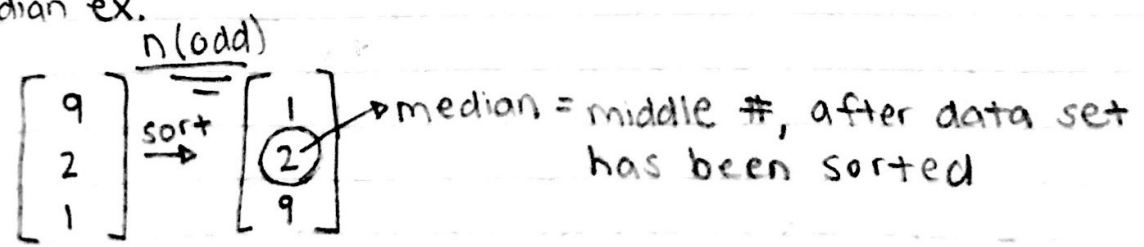
Bring reader and lecture notes to class

-the mean is the balance point

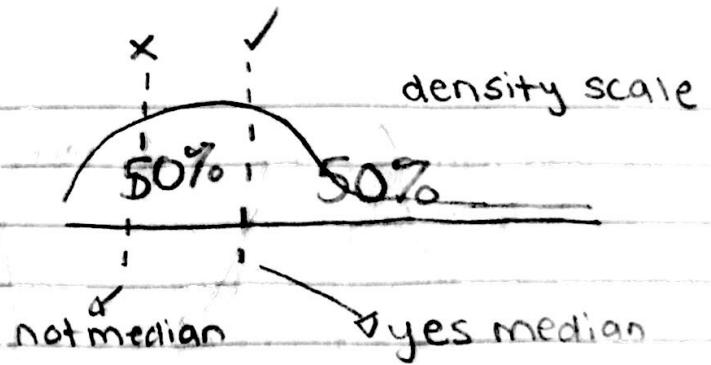
ex. L15 (butterfly winglength)



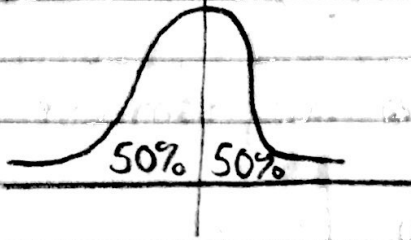
median ex.



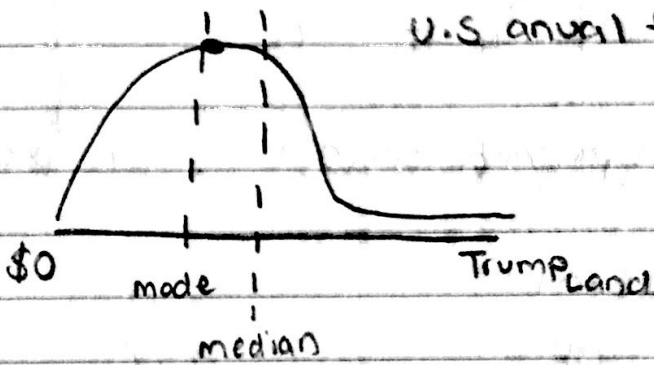
median - 50/50th point in data



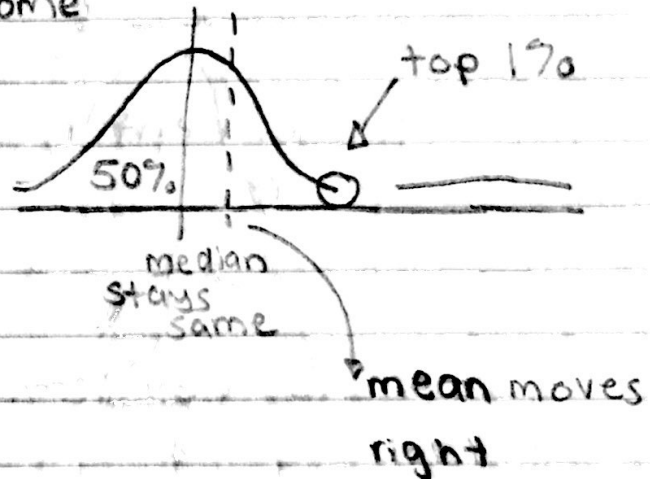
Symmetric unimodal



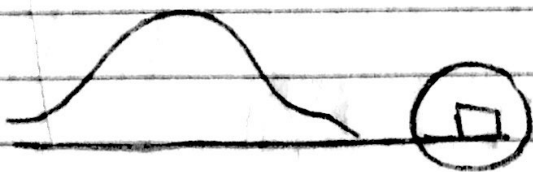
mode = point of symmetry
= mean
= median



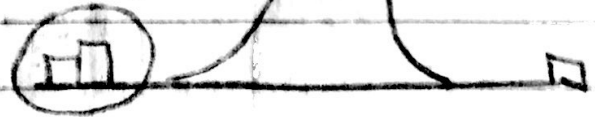
U.S. annual family income



right-tail outlier



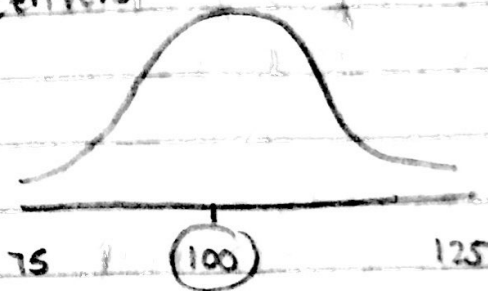
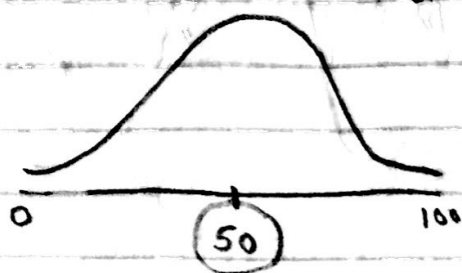
left-tail outliers



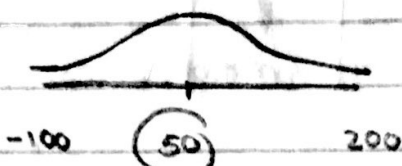
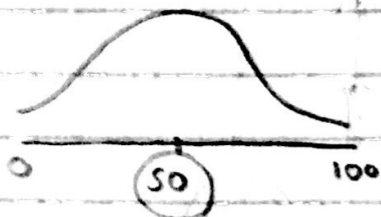
median is insensitive to outliers, but mean can be highly sensitive to outliers

* Look up page L-23-24 for more information
Ex. L-25 (Life expectancy of birds in captivity)

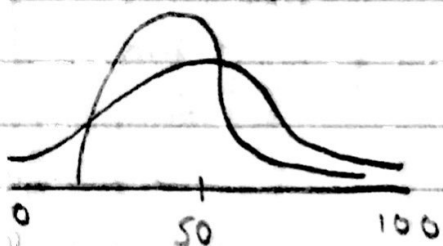
different centers



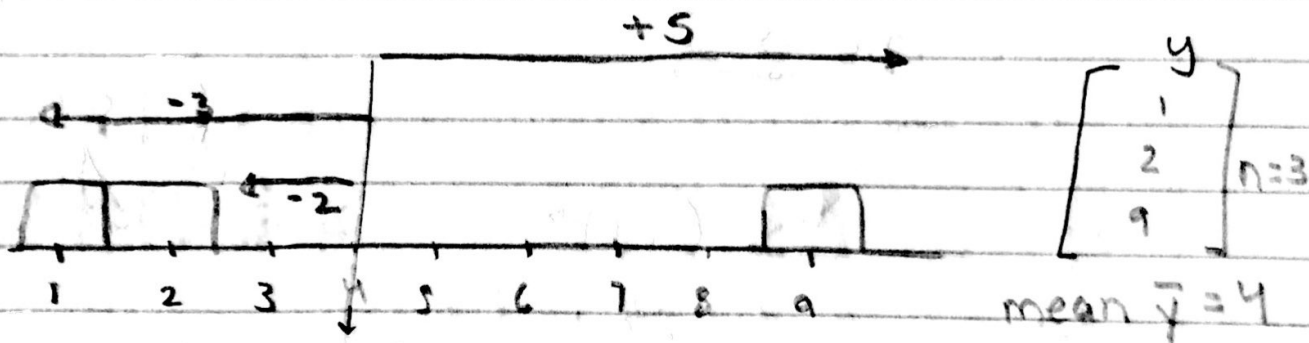
same spread, and basically same shape



same center, different spread, same basic shape



- same center
- same spread
- different shape



mean $\bar{y} = 4$

$$\begin{bmatrix} 1 \\ 2 \\ 9 \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

subtract $\bar{y} = 4$ from all data values

$$\begin{bmatrix} -3 \\ -2 \\ +5 \end{bmatrix} = \begin{bmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{bmatrix}$$

mean = 0

to avoid $\oplus + \ominus$ cancellation,

1) absolute values

derivatives from the means

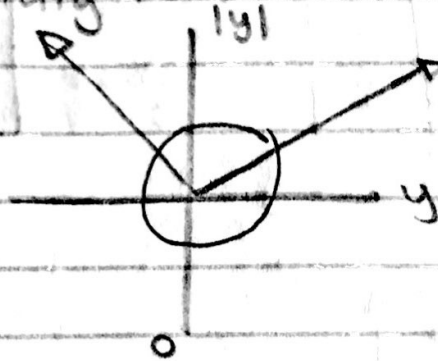
$$\begin{bmatrix} |-3| \\ |-2| \\ |5| \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} \quad \begin{bmatrix} |y_1 - \bar{y}| \\ \vdots \\ |y_n - \bar{y}| \end{bmatrix}$$

mad = mean absolute
de

mean 3.3 mean

MAD (sir arthur Eddington ~ 1910)

↳ not used
much today



$|y|$ is not
differentiable
at 0

$$\begin{bmatrix} \$1 \\ \$2 \\ \$9 \end{bmatrix} \xrightarrow[\text{mean } \$4]{\text{subtract } 4} \begin{bmatrix} -3 \\ -2 \\ +5 \end{bmatrix} \xrightarrow{\text{square}} \begin{bmatrix} (-3)^2 = 9\$ \\ (-2)^2 = 4\$ \\ (5)^2 = 25\$ \end{bmatrix}$$

mean = 12.7 $\2

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \xrightarrow[\text{mean } \bar{y}]{\text{subtract } \bar{y}} \begin{bmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{bmatrix} \xrightarrow{\text{square}} \begin{bmatrix} (y_1 - \bar{y})^2 \\ \vdots \\ (y_n - \bar{y})^2 \end{bmatrix}$$

$$(-3)^2 + (-2)^2 + (+5)^2 = (\text{sample}) \text{ variance } (s^2)$$

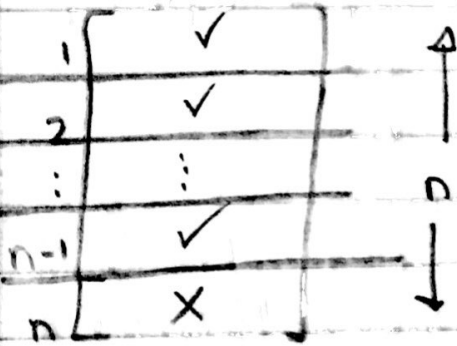
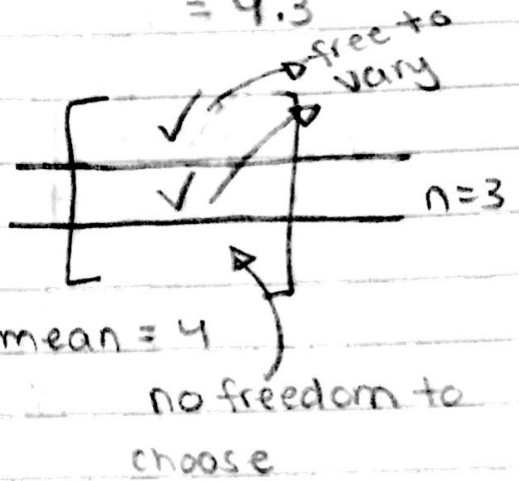
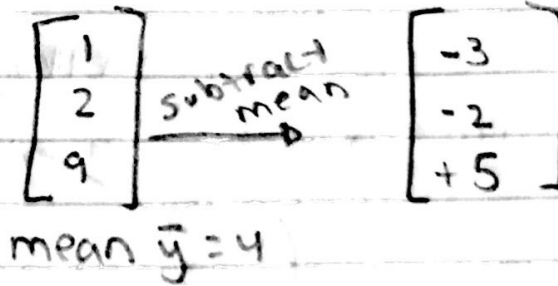
$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

(sample)
 Standard deviation = $\sqrt{\frac{\text{(sample variance)}}{n}} = s$
 (SD)

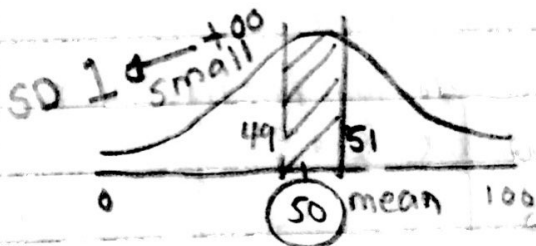
$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}$$

here $s = \sqrt{\frac{38}{2}} = 4.3$

why (n-1)
 not n!

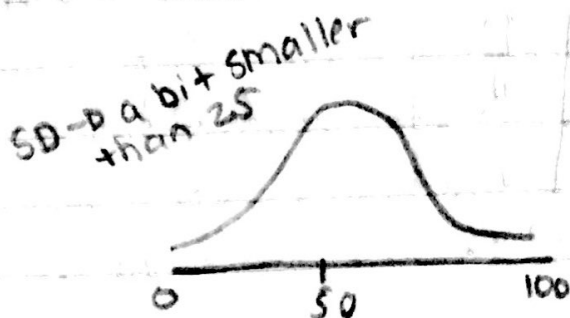
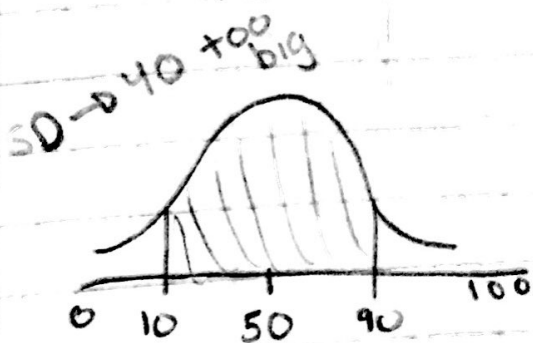


* a data set with n obs.
 only has (n-1) degrees of freedom for measuring spread



graphical interpretation of SD
 * empirical rule

- for virtually any data set, if you start at the mean + so (1) SD either way, you will capture about 2/3 (68%) most 95% (95%) almost all (99.7%) of the data in that interval



April 17, 2017

This time: normal curve
Next time: experimental design

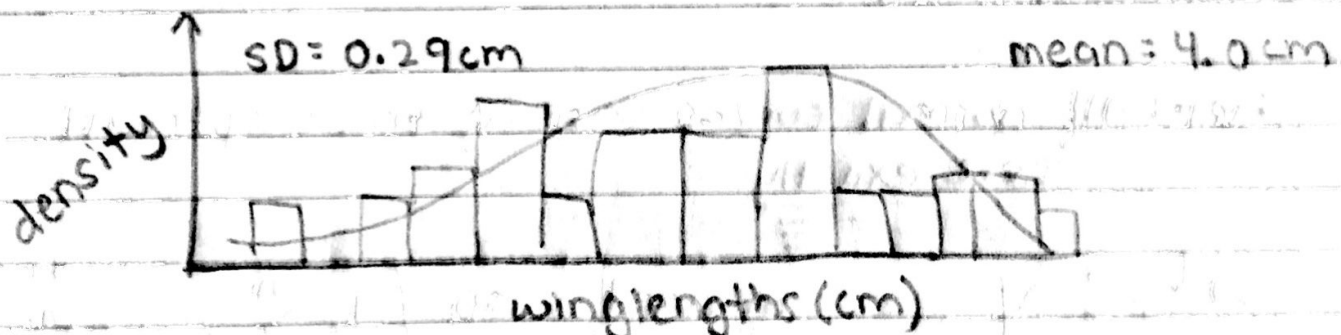
read: DD ch. 1-3(A)
ch. 1-6(B)

read: LN pp. 1-94

New HW due date: Friday 21 April @ 11:59pm

Get course materials packet soon if you don't already have it: will need it for discussion section week two

SG office hours: wed. 3-5pm SG BE 119, not 319



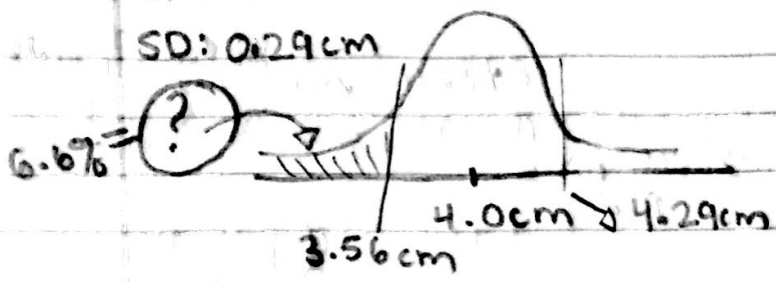
Q: what % of butterflies in data set had wing lengths ≤ 3.56 cm?

A₁: (exact) relative freq. from data: $\frac{2}{24} = 0.083$ or 8.3%

A₂: (exact) area under histogram < 3.56 cm = $\frac{2}{24} = 8.3\%$

Quetelet (1796 - 1874)

A₃: (approximate) work out area under bell curve (normal curve) that best approximates histogram to left of 3.56 cm



raw units for y is above average

↓

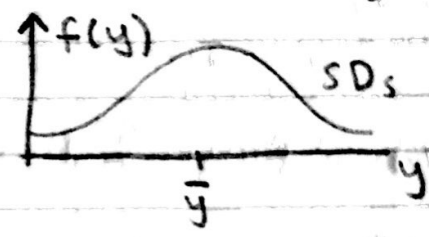
standard units =

(?) - 1 0 1

↓ -1.517

= -1.52

normal curve density formula

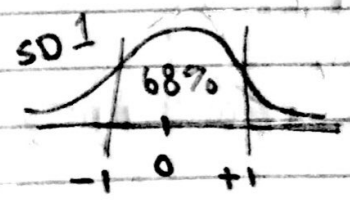


$$f(y) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2\sigma^2}(y-\bar{y})^2\right]$$

not integrable

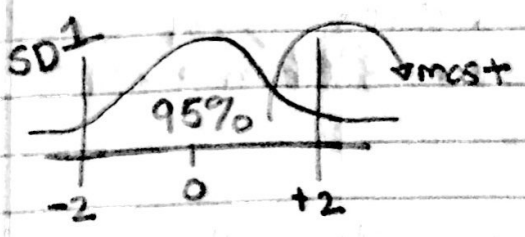
Isaac Newton = numerical integration
Gottfried Leibniz

Fact: All normal curves satisfy the Empirical Rule exactly



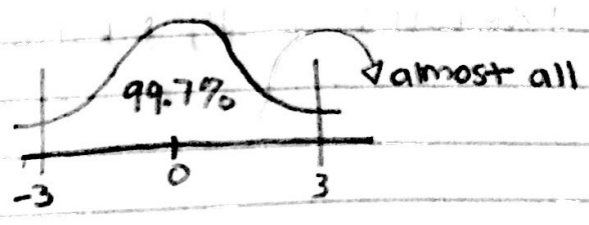
standard normal curve

$$SD = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}$$



SD = 0 ?
-negative SDs impossible

$$\begin{bmatrix} 4 \\ \vdots \\ 4 \end{bmatrix} \quad SD: 0$$



$$0 \leq SD < \infty$$

Look up pages L-34 through L-35

Facts about normal curve density scale:

- ① it's symmetric
- ② total area under curve = 100%

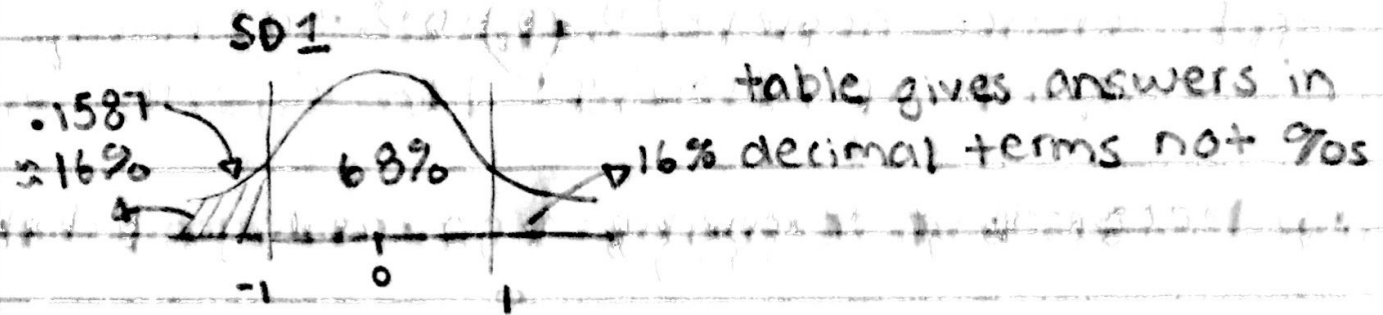
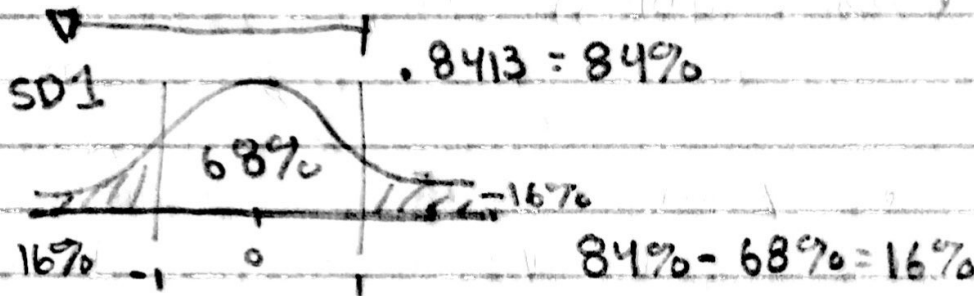


table gives answers in 16% decimal terms not %os



Converting from raw units to standard units

$$su = \frac{\# - \text{mean}}{SD}$$

$$z = \frac{y - \bar{y}}{s}$$

z scores (su) are unitless

$$\frac{3.56 - 4.00}{0.29} = \frac{-0.44}{0.29} = -1.517$$

area under normal curve

normal approx.

= normal

rather crude

approximation

to actual

6.6%

rel. freq.

8.3%

Convert y from standard to raw units

$$y = \bar{y} + s \cdot z$$

formula sheet:

Reader - 22 → 28

(L-69) (R-41)

read: science article (R-41) → (R-50) now, +
again in week 10

to decrease: uncertainty → get more good data