

this measure of
time: center & spread

read: JJ
book A ch. 1, 2
book B = ch. 1-3

AMS 7
9 OCT 18

next
time: normal
curve

lecture notes:
LN pp. ① - ②②

today
LN
pp. ②② +

new course deadline for HW I: 11:59 pm
next wed (8 days from now)

3 possible vertical scales for histograms

- ① raw-frequency: plot the counts
- ② relative-frequency: plot the %
- ③ density scale: when hist. are

plotted on density scale:

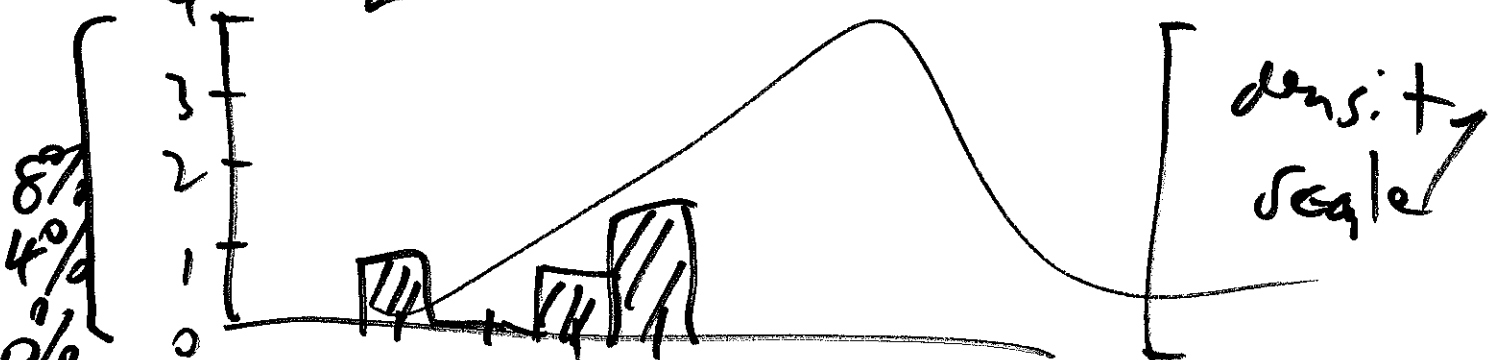
- (a) rel. freq \rightarrow area of histogram bars
- (b) total area under hist. = 100%

butterfly data

Convention:
all hist. from
now on are
implicitly on
density scale

value	row freq.	% (relative freq.)
3.3	1	$1/24 = 4\%$
3.4	0	0%
3.5	1	4%
3.6	2	8% ($\frac{2}{24}$)
...
4.5	1	4%
$n = 24$		100%

relative
freq. \downarrow
row
freq.



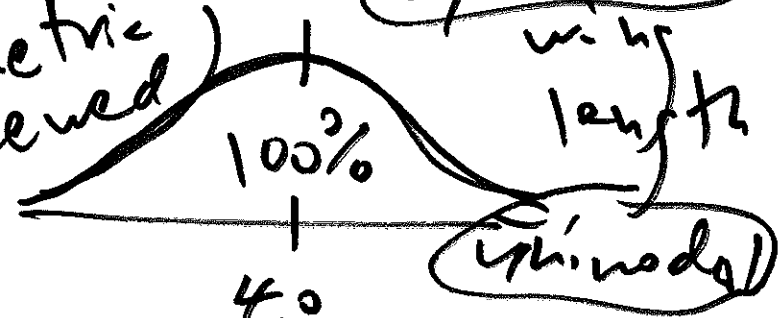
density
scale

positively
skewed



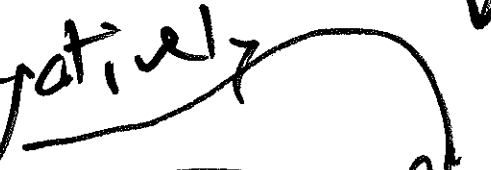
not
symmetric
(skewed)

symmetric
with
length



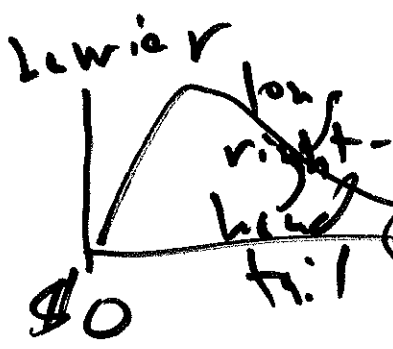
thin model

negatively
skewed

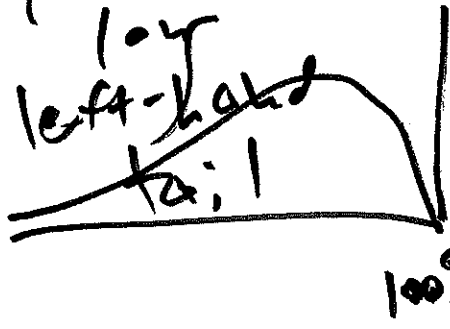
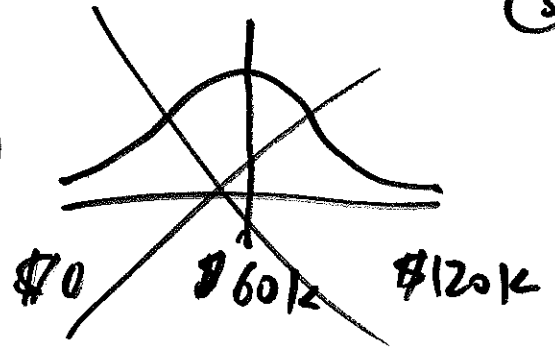


100%
|
4.0
cm

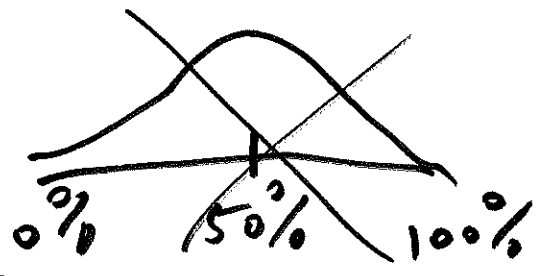
point
of
symmetry



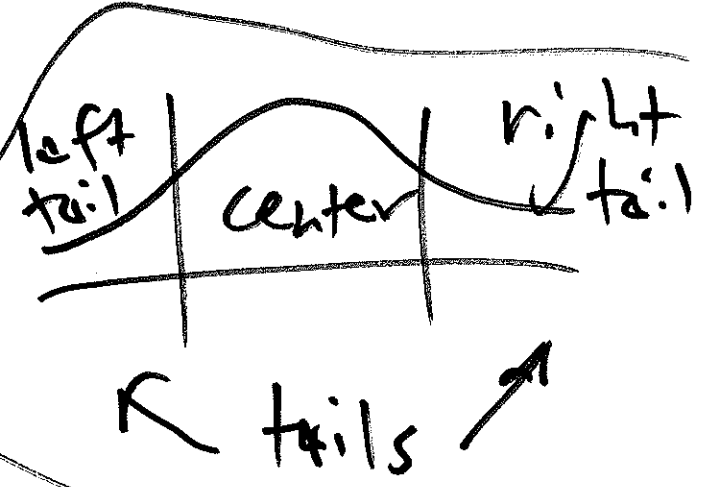
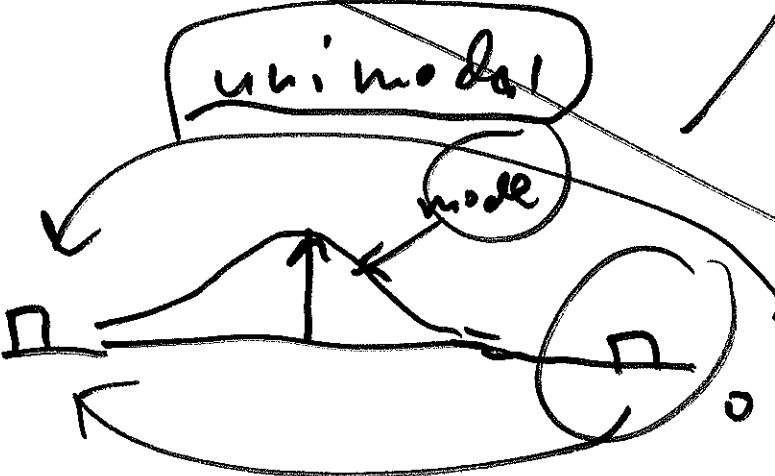
us. family income in 2017



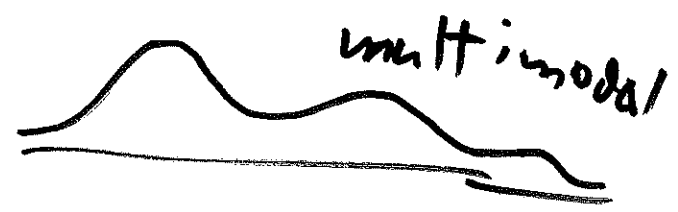
widtn score (%)



barrier

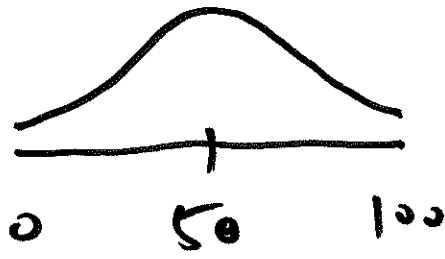


bimodal

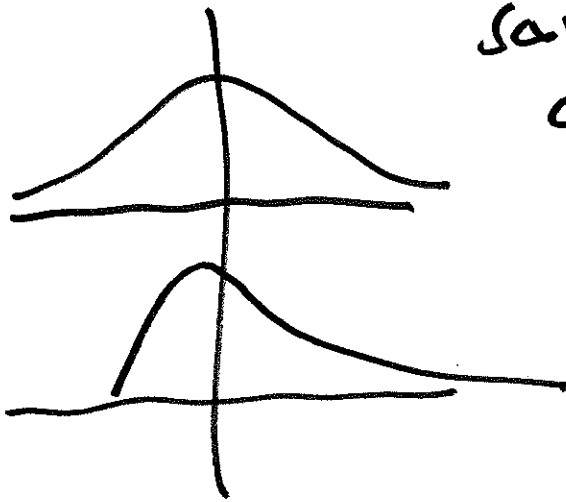
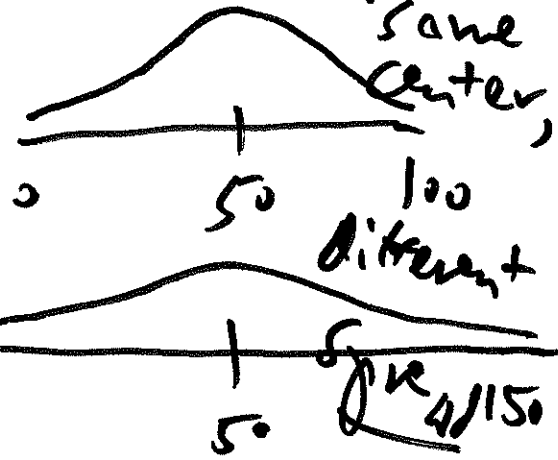
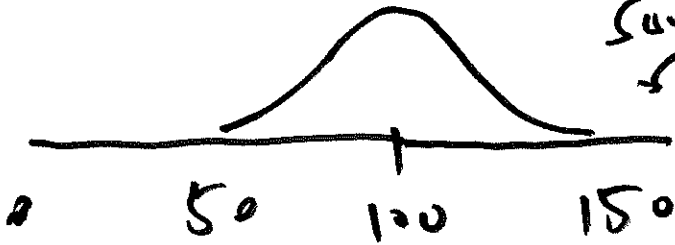


multimodal

same shape, (4)



different center,
same shape,
same spread



same center,
same spread,
different

qualitatively

shape

measures of center

L-15

quant. cont. ratio.

y_i	length (cm)
y_1	4.4
y_2	3.6
\vdots	\vdots
y_k	3.9

$n = 24$

(1) mean (average)

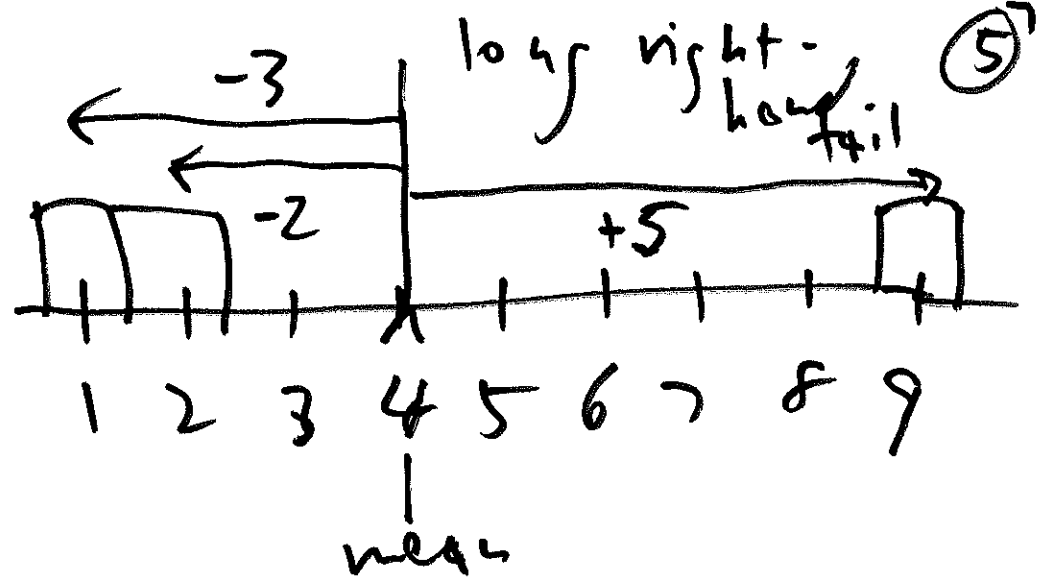
(2)

(3) mode

$$\text{mean } \bar{y} = \frac{4.4 + \dots + 3.9}{24} = 4.0 \text{ cm}$$

$$\begin{matrix} y_1 \\ \vdots \\ y_n \end{matrix} \begin{bmatrix} 1 \\ 2 \\ 9 \end{bmatrix} \quad n=3$$

mean $\bar{y} = 4$



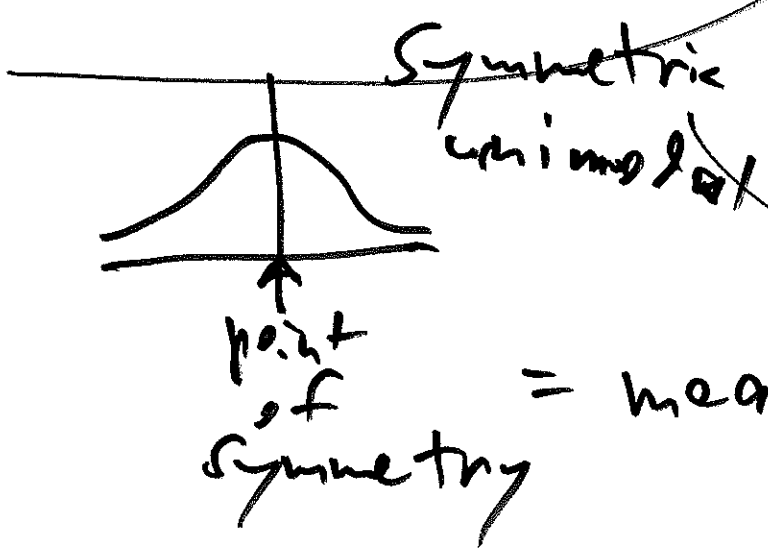
$$\begin{bmatrix} 1 \\ 2 \\ 9 \\ \text{mean } 4 \end{bmatrix} \xrightarrow[\text{mean } 4]{\text{subtract}} \begin{bmatrix} -3 \\ -2 \\ +5 \\ \text{mean } 0 \end{bmatrix}$$

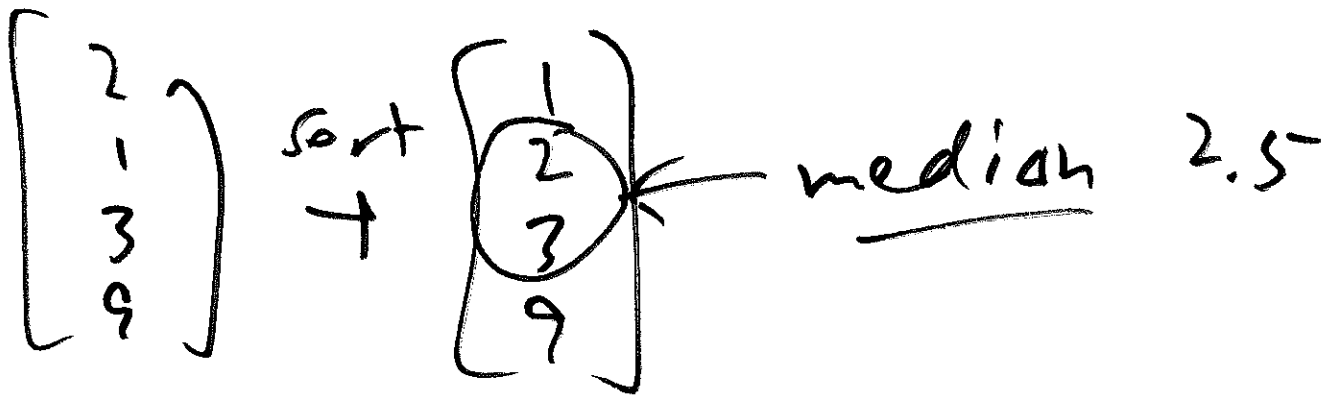
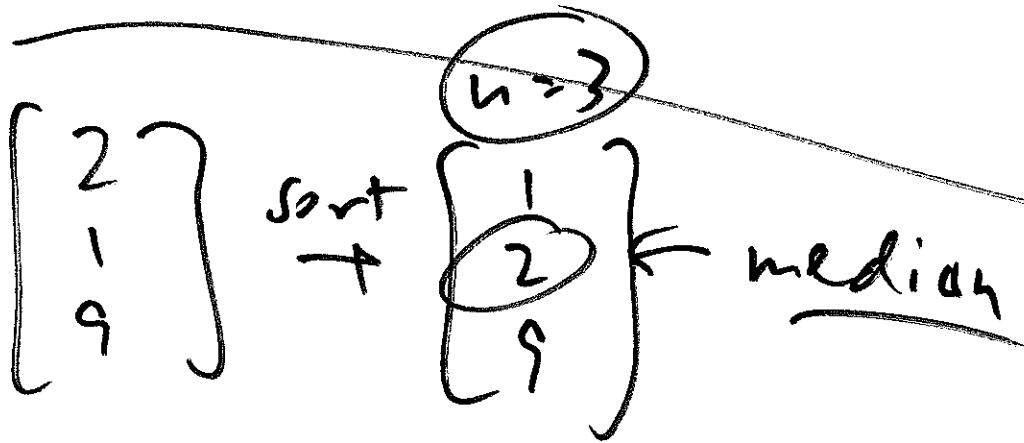
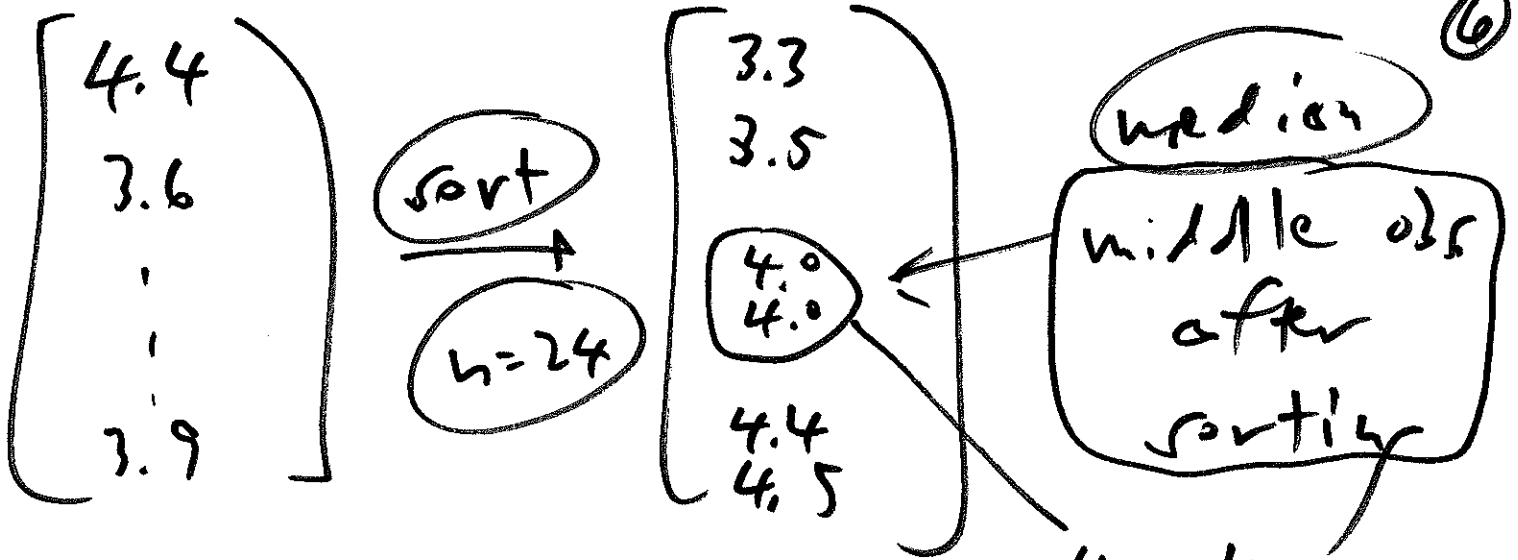
deviations from the mean

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \\ \text{mean } \bar{y} \end{bmatrix} \xrightarrow[\text{mean } \bar{y}]{\text{subtract}} \begin{bmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \\ \text{mean } 0 \end{bmatrix}$$

graphical interpretation of mean:

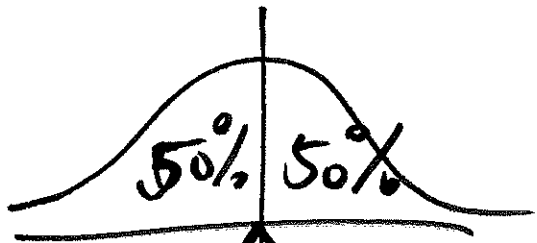
center of gravity
= balance point





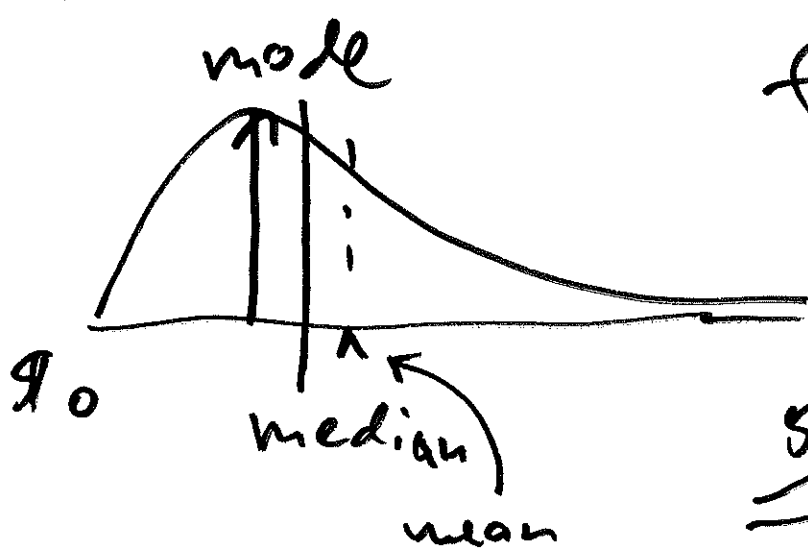
graphical interpretation of median

50% / 50% point in data in relative frequency terms

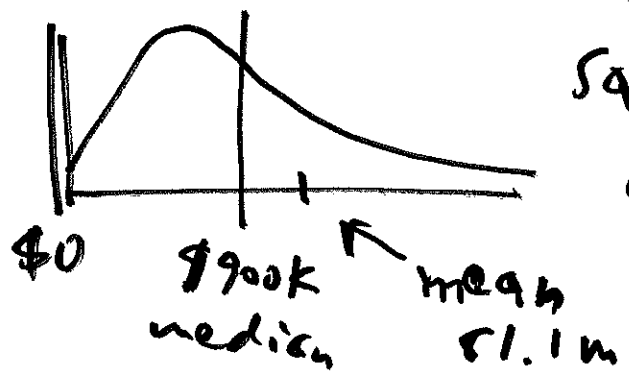
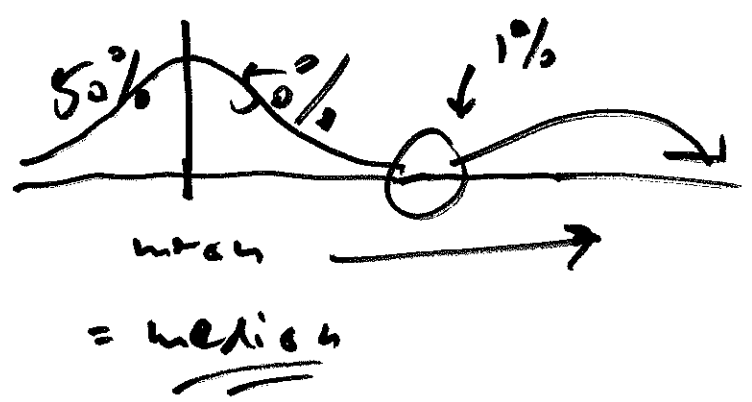


point of symmetry

= mean = mode = median



family income



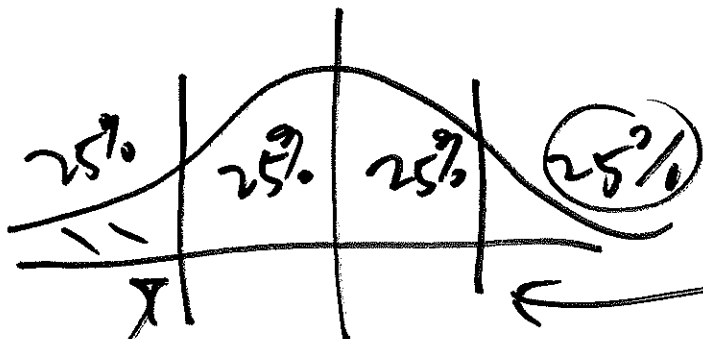
sale price of S.C. houses

lower case sigma

upper case sigma

$$\bar{y} = \frac{y_1 + \dots + y_n}{n} = \frac{1}{n} \sum_{i=1}^n y_i$$

index of summation



~~75th~~ 75th percentile = 0.75 quantile

median = 50th percentile = 0.5 quantile

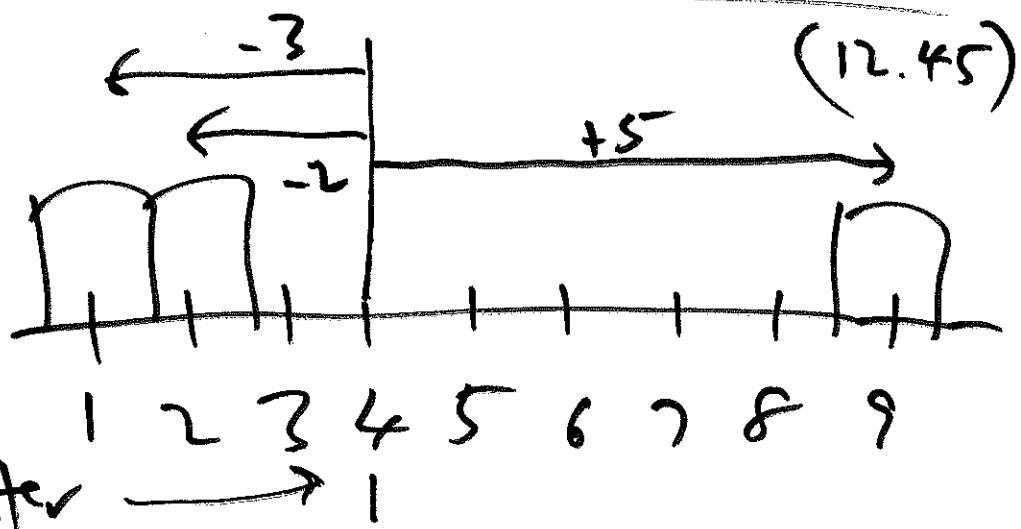
25th

percentile = 0.25 quantile

influence of outliers on mean:

mean is pulled by the tail

measures of spread
 typical amount by which each # differs from center



sample

$$\begin{bmatrix} 1 \\ 2 \\ 9 \end{bmatrix}$$

mean 4

subtract

$$\begin{bmatrix} -3 \\ -2 \\ +5 \end{bmatrix}$$

mean 0

absolute value

$$\begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix}$$

mean $\frac{10}{3} = 3.3$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

mean \bar{y}

subtract

$$\begin{bmatrix} y_1 - \bar{y} \\ y_2 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{bmatrix}$$

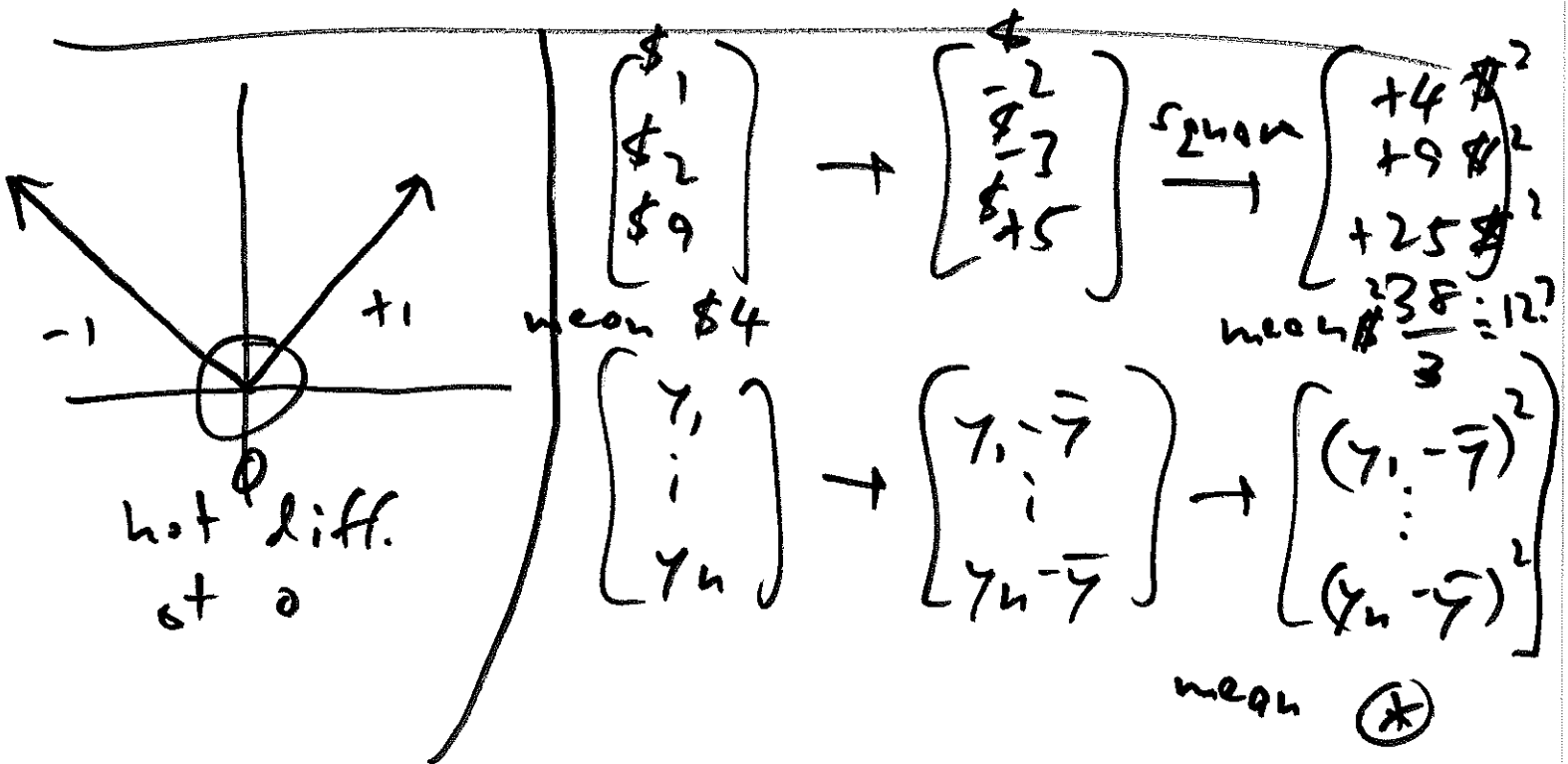
mean 0

abs. value

$$\begin{bmatrix} |y_1 - \bar{y}| \\ \vdots \\ |y_n - \bar{y}| \end{bmatrix}$$

mean

$$\frac{1}{n} \sum_{i=1}^n |y_i - \bar{y}| = (\text{MAD}) \text{ mean absolute deviation}$$

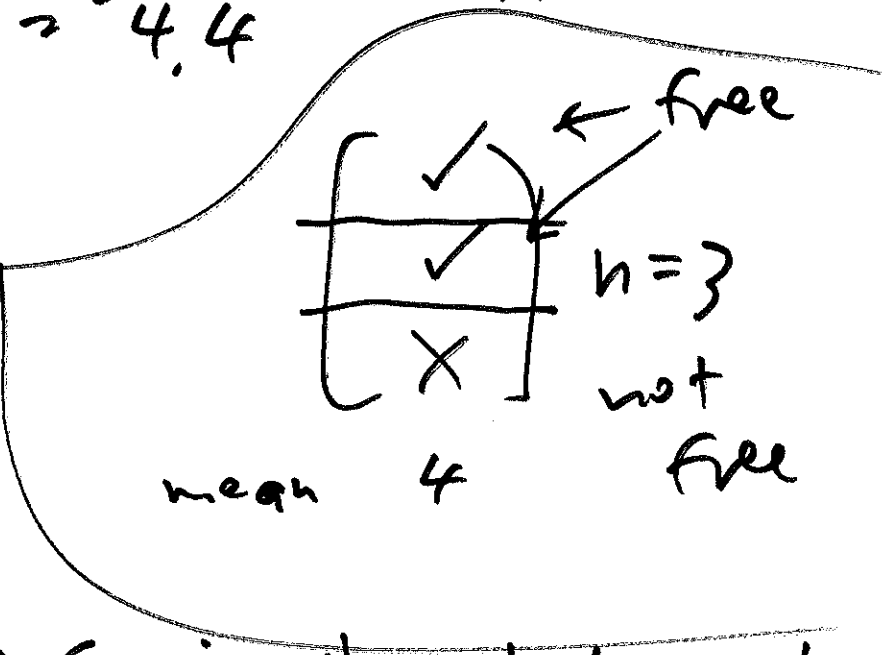


⊗ $s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = \text{(sample) variance}$

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

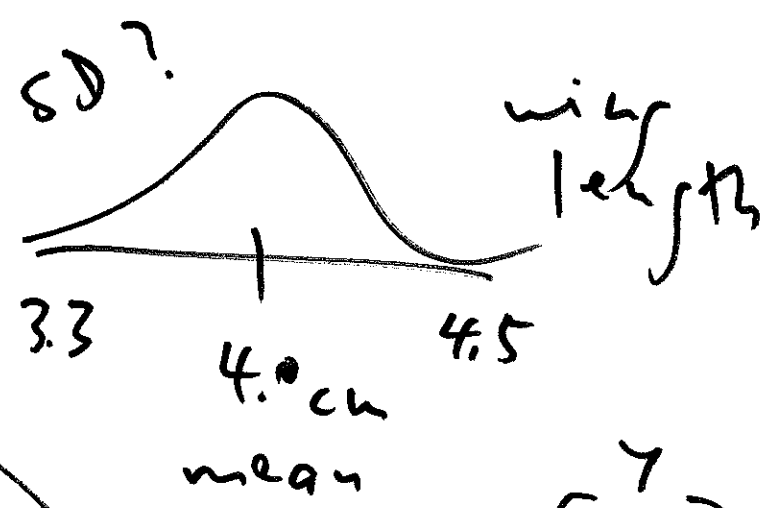
SD = $\sqrt{\frac{38}{2}} = 4.4$ (100 yrs old)

The data set has $n=3$



observations in it, but only $(n-1) = 2$ degrees of freedom for measuring spread

graphical interpretation of SD



empirical rule

start at mean,

go $\left(\begin{matrix} 1 \\ 2 \\ 3 \end{matrix} \right)$ SD (s) either way:

you will capture

of the data

- $n=24$
- 7
 - 4.4
 - 3.6
 - ...
 - 3.9

mean $\bar{y} = 4.0$ cm
SD $s = 0.29$

(about $\frac{2}{3}$) (68%)
most (95%)
almost all (99.7%)

0.5 too big (normal curve)
0.3 about right
0.1 too small