

10/9/18

# Wave Moretto

## STATS Lecture #4: Measures of Center + Spread

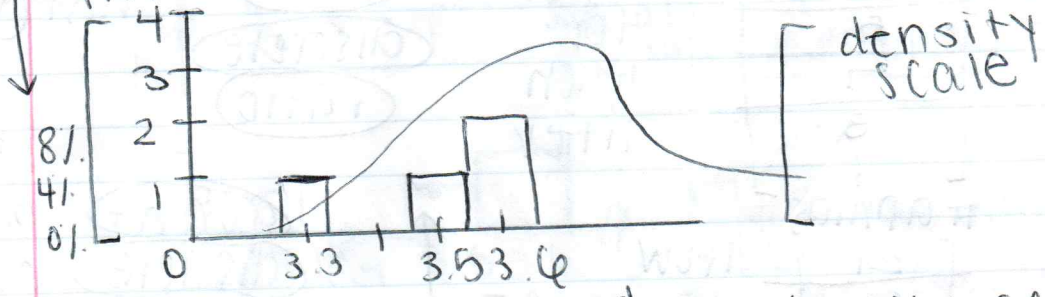
° 3 possible vertical scale histograms

- ① Raw-frequency: plot the counts
  - ② relative-frequency: plot the %.
  - ③ Density scale
- butterfly data

value	raw freq.	% (relative frequency)
3.3	1	1/24 = 4%
3.4	0	0%
3.5	1	4%
3.6	2	8%
4.5	1	4%
n = 24		100%

relative freq.

raw freq.



• When hist. are plotted on density scale:

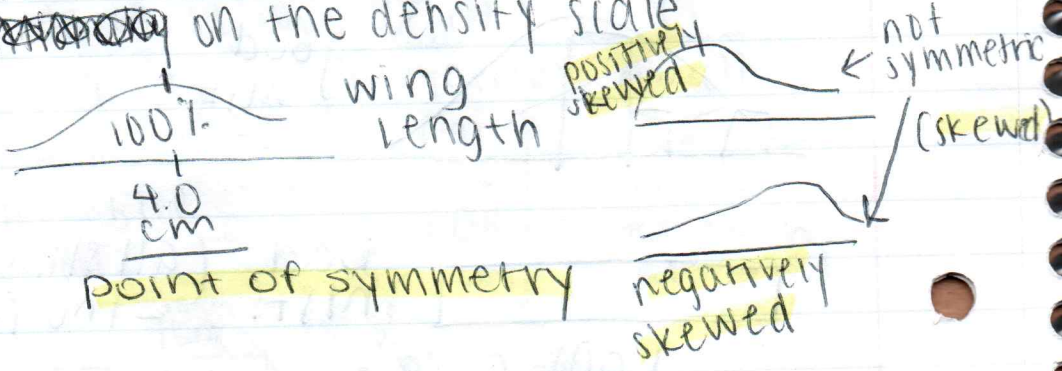
- a) rel. freq  $\leftrightarrow$  area of hist. bars (curve)
- b) total area under hist. = 100%

• Convention: all hist. from now on, are

implicitly

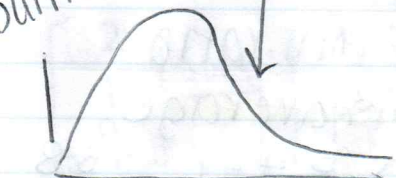
~~on the density scale~~ on the density scale

ex)



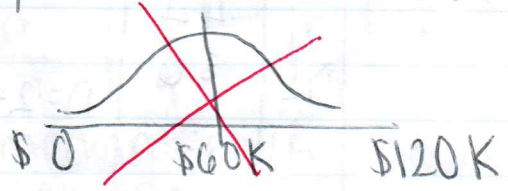
long right hand tail

barrier



U.S. family income in 2017

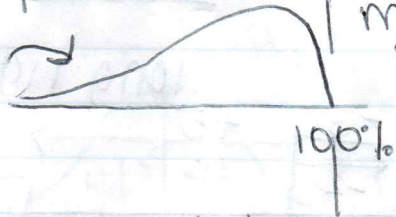
\$0 Bill Gates



\$0 \$60K \$120K

long left hand tail

midterm scores (%)



100%



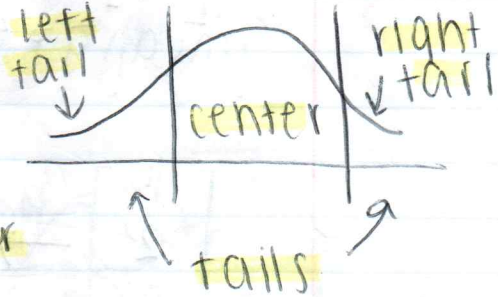
0% 50% 100%

unimodal

mode

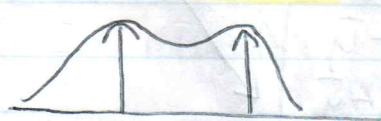


outlier

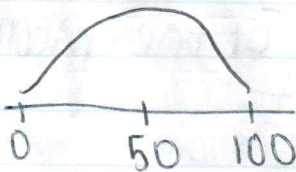
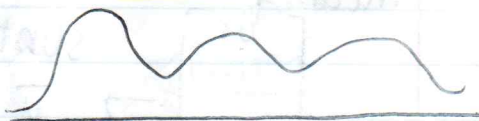


tails

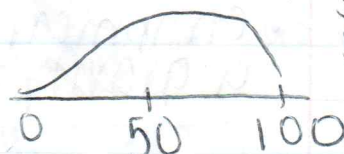
bimodal



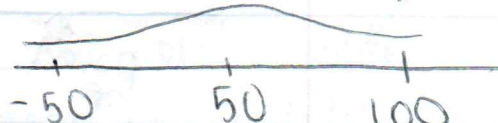
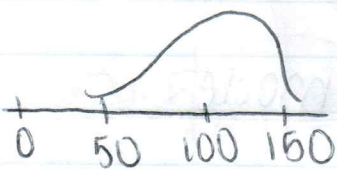
multimodal



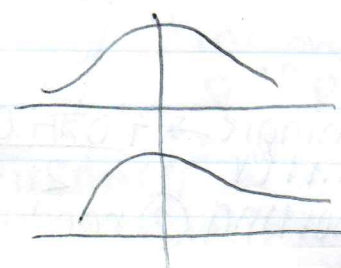
different center, same shape, same spread



same shape, same center, different spread



same center, same spread, different shape



qualitatively

Measures of center: (L-15)

$$y_1 \begin{bmatrix} 4.4 \\ 3.6 \\ \vdots \\ 3.8 \end{bmatrix} n=24$$

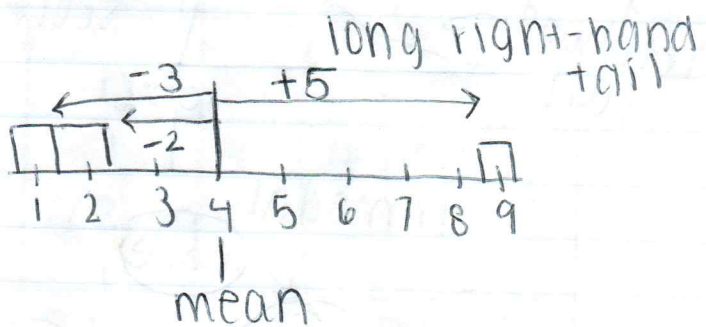
quant. cont. ratio

① mean/average

$$\text{mean } \bar{y} = \frac{4.4 + \dots + 3.8}{24} = \boxed{4.0 \text{ cm}}$$

$$\begin{bmatrix} 1 \\ 2 \\ 9 \end{bmatrix} n=3$$

mean  $\bar{y} = 4$



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

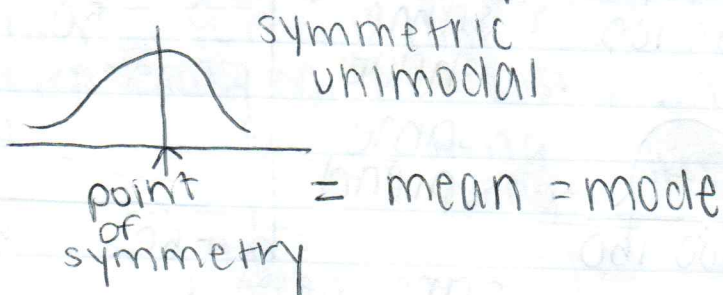
mean 4

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

mean 0

deviations from the mean

Graphical interpretation of the mean: center of gravity = balance point



$$\begin{bmatrix} 4.4 \\ 3.6 \\ \vdots \\ 3.9 \end{bmatrix} \xrightarrow{\text{sort}} \begin{bmatrix} 3.3 \\ 3.5 \\ 4.0 \\ 4.0 \\ 4.4 \\ 4.5 \end{bmatrix}$$

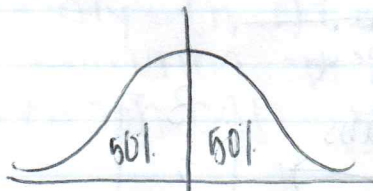
middle after sorting  $\rightarrow \frac{4.0 + 4.0}{2} = 4.0$

② median

$n=3$

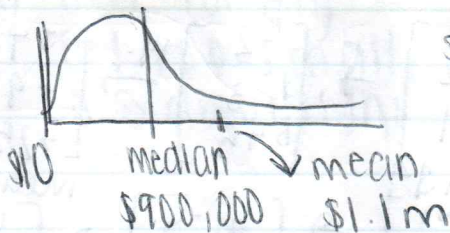
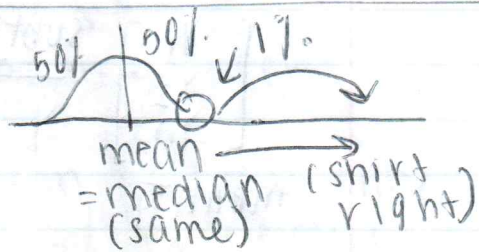
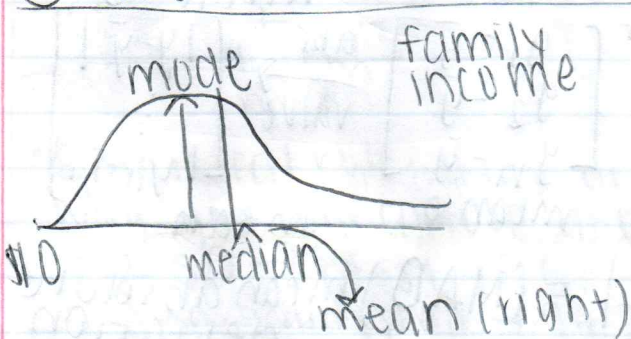
$$\begin{bmatrix} 2 \\ 1 \\ 9 \end{bmatrix} \xrightarrow{\text{sort}} \begin{bmatrix} 1 \\ \textcircled{2} \\ 9 \end{bmatrix} \leftarrow \text{median} \quad \begin{bmatrix} 2 \\ 1 \\ 3 \\ 9 \end{bmatrix} \xrightarrow{\text{sort}} \begin{bmatrix} 1 \\ \textcircled{2} \\ \textcircled{3} \\ 9 \end{bmatrix} \rightarrow \text{median} = 2.5$$

Graphical interpretation of median: 50/50 point in data in relative frequency terms



point of symmetry = mean = mode = median

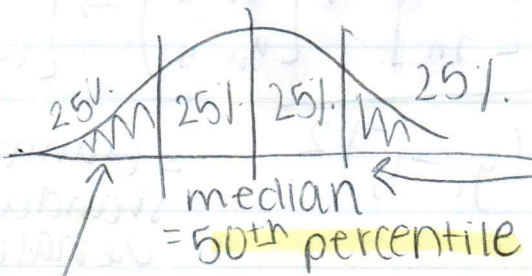
### ③ Mode



sale price of SC housing

$\sigma$  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 percentile = 0.75 quantile

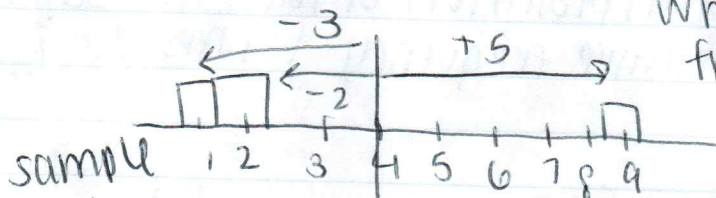
25th percentile = 0.25 quantile

• Influence of outliers on the mean

• Mean is pulled by the tail

Measures of spread:

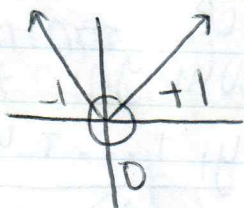
- typical amount by which each # differs from center



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

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \xrightarrow[\text{mean } \bar{y}]{\text{subtract}} \begin{bmatrix} y_1 - \bar{y} \\ y_2 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{bmatrix} \xrightarrow[\text{mean}]{\text{abs value}} \begin{bmatrix} |y_1 - \bar{y}| \\ |y_2 - \bar{y}| \\ \vdots \\ |y_n - \bar{y}| \end{bmatrix}$$

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



not diff at 0

$$\begin{bmatrix} 1 \\ 2 \\ 9 \end{bmatrix} \xrightarrow[\text{mean } 4]{\text{subtract}} \begin{bmatrix} -2 \\ -3 \\ +5 \end{bmatrix} \xrightarrow[\text{mean: } 12.7]{\text{sq.}} \begin{bmatrix} +4 \\ +9 \\ +25 \end{bmatrix}$$

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

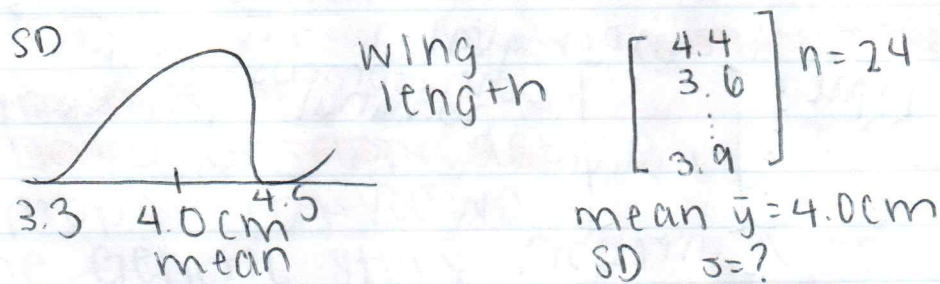
$$\text{mean: } s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2} = (\text{sample standard deviation (SD)})$$

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

$$SD = \sqrt{\frac{1^2 38}{2}} = 4.4 \quad \begin{matrix} \left[ \begin{array}{c} \checkmark \\ \checkmark \\ \bar{x} \end{array} \right] \begin{matrix} \leftarrow \text{free} \\ n=3 \\ \leftarrow \text{not free} \end{matrix} \\ \text{mean} \end{matrix}$$

(100 yrs old)

- The data set has  $n=3$  observations in it, but only  $(n-1)=2$  degrees of freedom for measuring spread
- Empirical interpretation of SD:



- **Empirical rule:** start at mean go 1 SD either way: you will capture about  $\frac{2}{3}$  of the data
- 2 SD  $\rightarrow$  most  $\rightarrow$  95%
- 3 SD  $\rightarrow$  almost all  $\rightarrow$  99.7%. (from normal curve)

ex) 0.5 too big  
0.1 too small  
0.3 about right