

this time: standard deviation; ^(SD) the normal curve

lead: PP ch. 1-3 (A) ^{latest} 1-6 (B) ^{earlier}

AMS 7 14 Apr 17

next time: controlled experiments & observational studies

LN pp. 1-94

ecomm 5 is out;

our official note-taker is Kaitlyn Abercrombie

Canvas. uesc.edu is where you'll submit homework

median

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

median = middle #, after data set has been sorted

quant. \uparrow $\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$ \uparrow

mean \bar{y}

median \tilde{y}

$\begin{bmatrix} 9 \\ 1 \\ 4 \\ 2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 2 \\ 4 \\ 9 \end{bmatrix}$ $n=4$

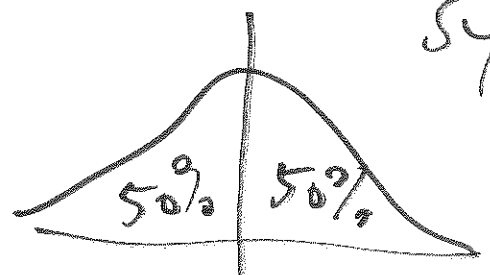
median = mean of middle 2 #s (n even)

graphical interpretation of median

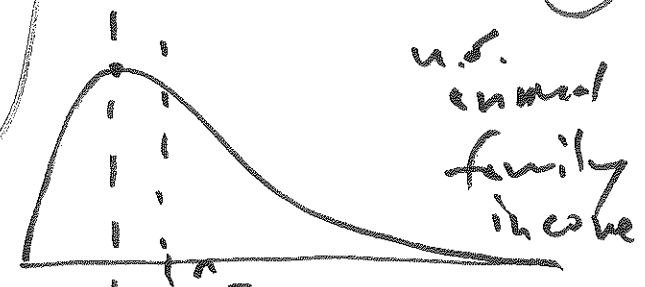
50%/50% point in data

(Density Scale)

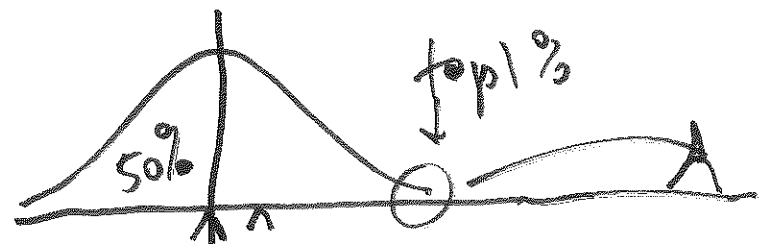
Symmetric unimodal



mode =
point of
symmetry
= mean
= median



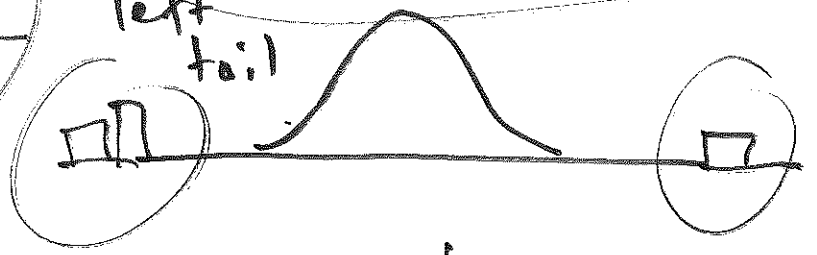
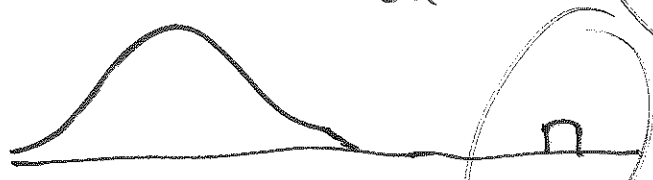
mode
median
mean



median stays same
mean moves right

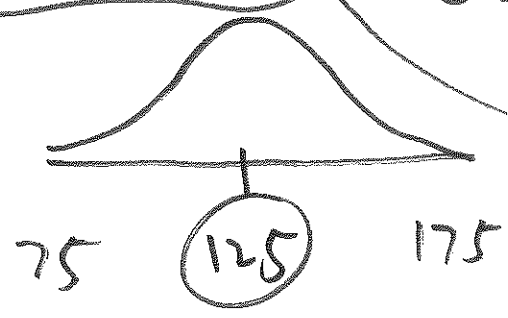
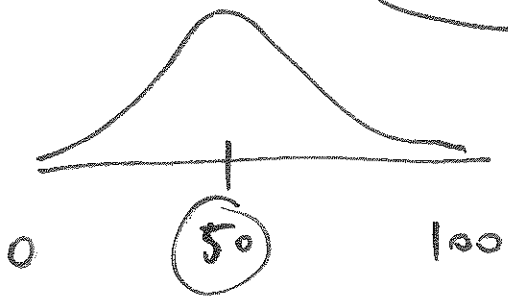
right tail outlier

left tail



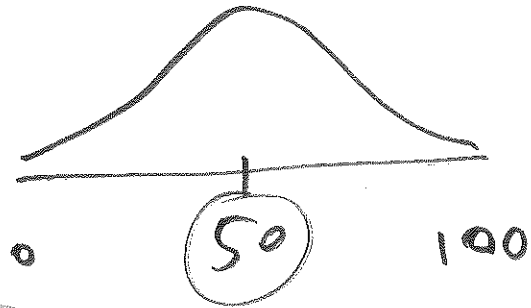
median insensitive to outliers, but mean can be highly sensitive to outliers

(different center)



same spread

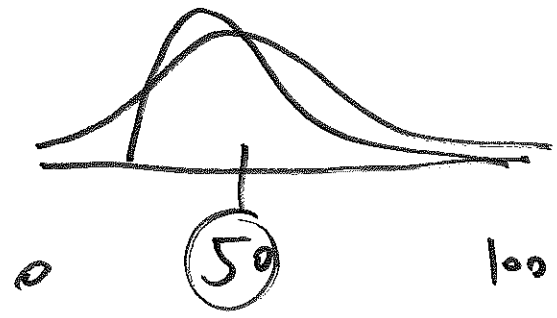
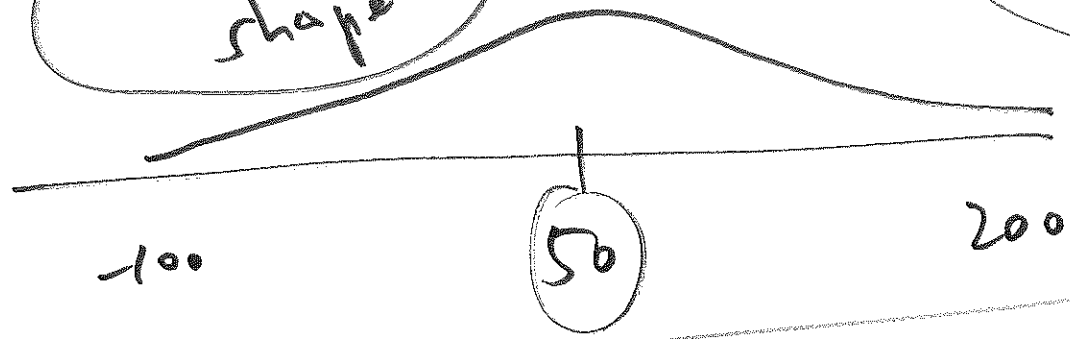
same shape



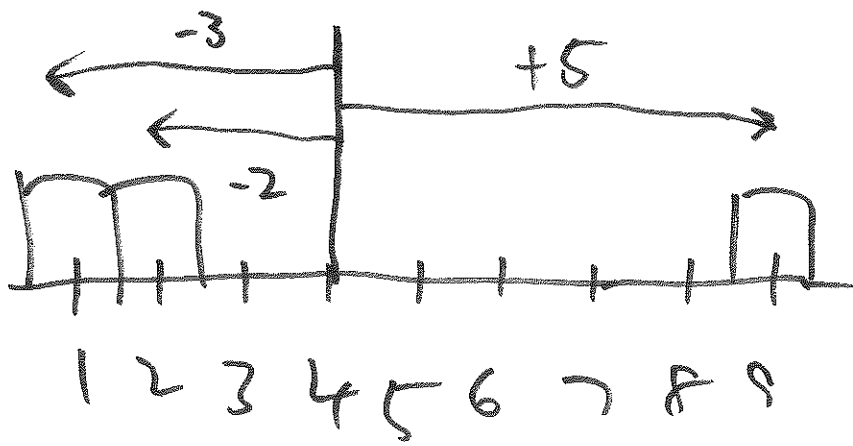
same center

same shape

different spread



same center
same spread
different shape



$$\begin{bmatrix} 1 \\ 2 \\ 9 \end{bmatrix} \quad n=3 \quad (4)$$

mean $\bar{y} = 4$

mean
 $\bar{y} = 4$

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

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

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

mean 0

deviations from the mean

to avoid \oplus & \ominus cancellation,

absolute values

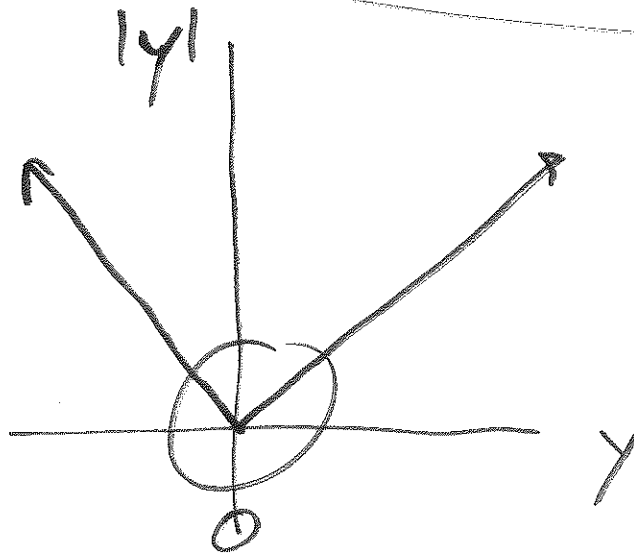
$$\begin{bmatrix} |1-3| \\ |1-2| \\ |+5| \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} \begin{bmatrix} |y_1 - \bar{y}| \\ \vdots \\ |y_n - \bar{y}| \end{bmatrix}$$

MAD = mean absolute deviation

mean 3.3

MAD (Sir Arthur Eddington) (5) ~1910

not used much today



$|y|$ is not differentiable at 0

$$\begin{array}{l} \left[\begin{array}{c} \$1 \\ \$2 \\ \$9 \\ \$4 \end{array} \right] \xrightarrow[\text{mean } 4]{\text{subtract}} \left[\begin{array}{c} \$-3 \\ \$-2 \\ \$5 \end{array} \right] \xrightarrow[\text{mean } 12.78^2]{\text{square}} \left[\begin{array}{c} (-3)^2 = +9 \\ (-2)^2 = +4 \\ (+5)^2 = +25 \end{array} \right] \end{array}$$

$$\begin{array}{l} \left[\begin{array}{c} y_1 \\ \vdots \\ y_n \end{array} \right] \xrightarrow[\text{mean } \bar{y}]{\text{subtract}} \left[\begin{array}{c} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{array} \right] \xrightarrow{\text{square}} \left[\begin{array}{c} (y_1 - \bar{y})^2 \\ \vdots \\ (y_n - \bar{y})^2 \end{array} \right] \end{array}$$

$$\frac{(-3)^2 + (-2)^2 + (+5)^2}{3-1} = (\text{sample}) \text{ variance} = 5^2$$

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

(sample) standard deviation = (SD)

(sample) variance = S

⑥

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

here

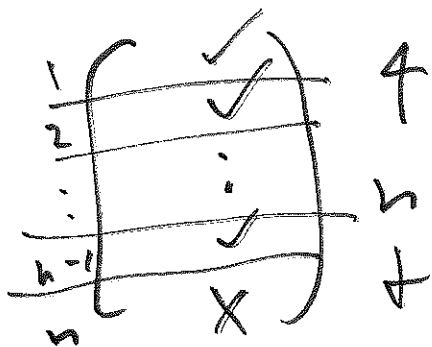
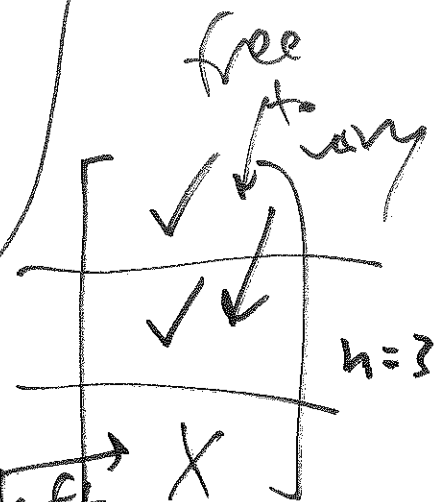
$$S = \sqrt{\frac{38}{2}} = 4.3$$

why (n-1) not n!

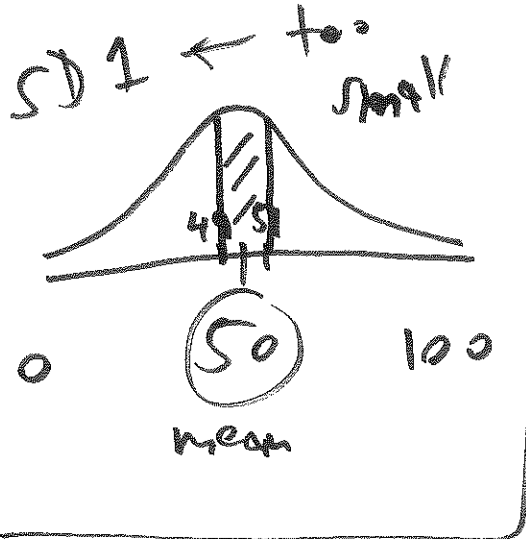
mean $\bar{y} = 4$

subtract \rightarrow mean

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



a dataset 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

& go $\begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}$ SD(s), either way, you

will usually capture $\begin{pmatrix} \text{about } 2/3 & 68\% \\ \text{most} & 95\% \\ \text{almost all} & 99.7\% \end{pmatrix}$ of the data in that interval

