

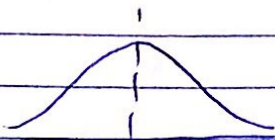
9<sup>th</sup> October 2018

Nikhila Kadiyala

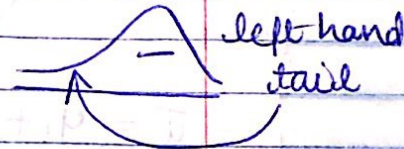
3 possible vertical scales for Histograms

value	raw freq	% (relative frequency)	density scale
3.3	1	4%	(a) relative frequency (b) total area under Histogram = 100%
3.4	0	0%	
3.5	1		
3.6	2	8%	

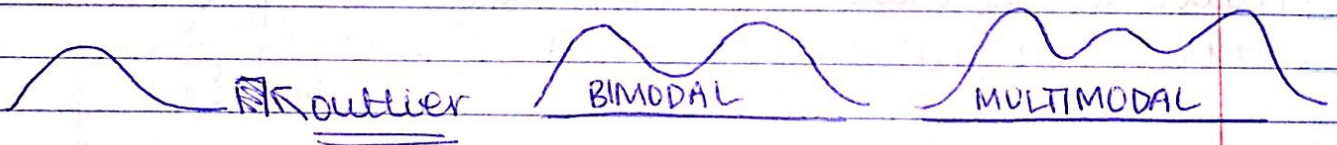
convention: all histograms are implicitly on the density scale



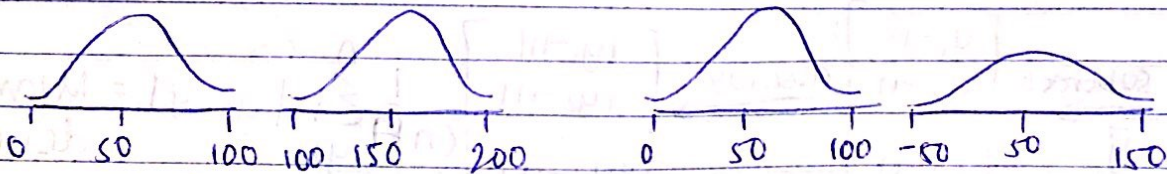
Not symmetric = skewed



• Skewness occurs when there is a barrier beyond/below which you can't go.



Bell Curve - symmetric, unimodal.



same shape, different center

same shape, same center  
different spread

Measure of center:

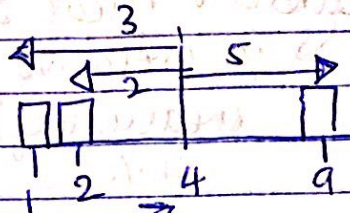
- (1) Mean (average)
- (2) Median (sort and the middle # is median)
- (3) Mode



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

$n=3$

Mean  $\bar{y} = 4$



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

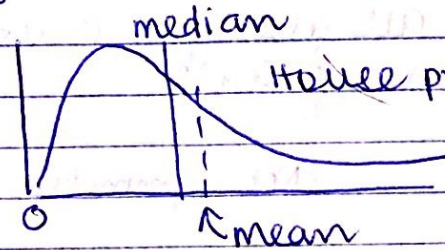
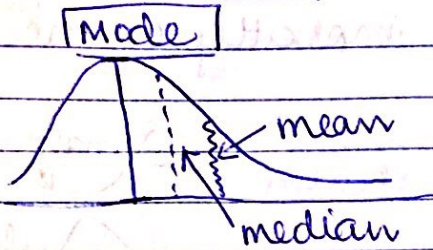
subtract  
4

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

↑  
deviations  
of the mean

mean of  $(y_i - \bar{y}) = 0$

center of gravity (balance point)

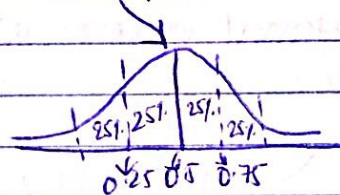


House prices in Sc.

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

Median: 50% percentile  
0.5 quantile

Influence of outliers on the mean  
→ Mean is pulled by the tail.



Measures of spread: Typical amount by which each # differs from centre.

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \end{bmatrix}$$

subtract  
 $\bar{y}$

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

absolute

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

$$\frac{1}{n} \sum_{i=1}^n |y_i - \bar{y}| = \text{Mean Absolute deviation}$$

↓ Square

$$\begin{bmatrix} (y_1 - \bar{y})^2 \\ (y_2 - \bar{y})^2 \\ \vdots \\ (y_n - \bar{y})^2 \end{bmatrix}$$

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

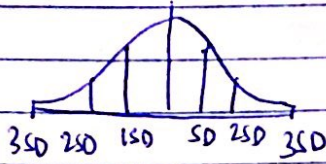
$$\sqrt{\text{Variance}} = \text{Standard deviation } (s)$$



• Graphical interpretation of SD.

EMPIRICAL RULE:

① start @ mean, go 1 SD either way,  
you will capture  $\frac{2}{3}$  of the data ( $\approx 68\%$ )



② start @ mean, go 2 SDs either way  
capture most data (95%)

③ start @ mean, go 3 SDs either way,  
capture almost all data (99.7%)