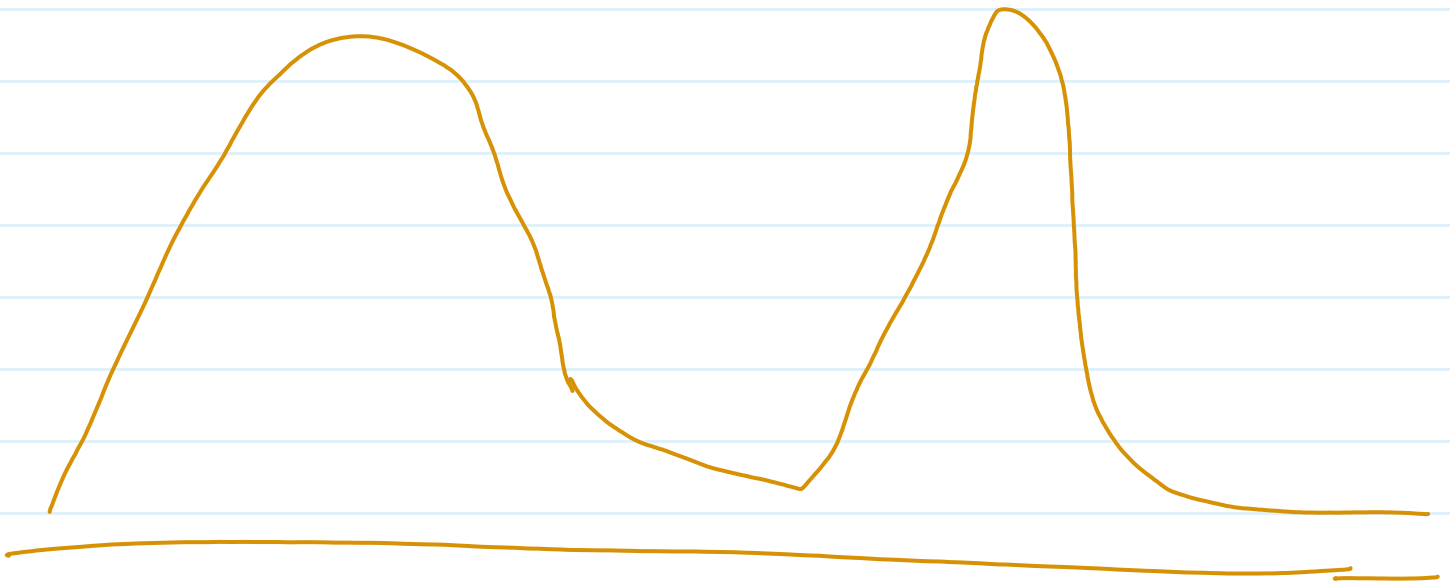


# \* Center limit theorem

The CLT states that regarding of the shape of population dist. the dist. of sample means will be approximately normal



$$n_1 - \bar{X}_1$$

$$n_2 - \bar{X}_2$$

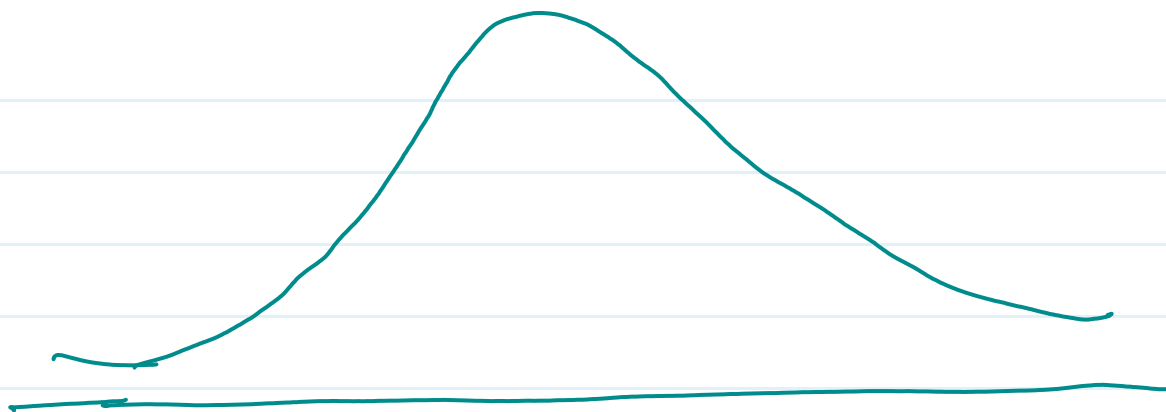
$$n_3 - \bar{X}_3$$

$$\vdots$$

$$n_{50} - \bar{X}_{50}$$

$$[\bar{X}_1, \bar{X}_2, \bar{X}_3, \dots, \bar{X}_{50}]$$

New  
sample  
data



g - population size - 10k

Sample size  $n \geq 30$

more the sample you will take  
more the plot will normally Dist.

## ★ Histogram

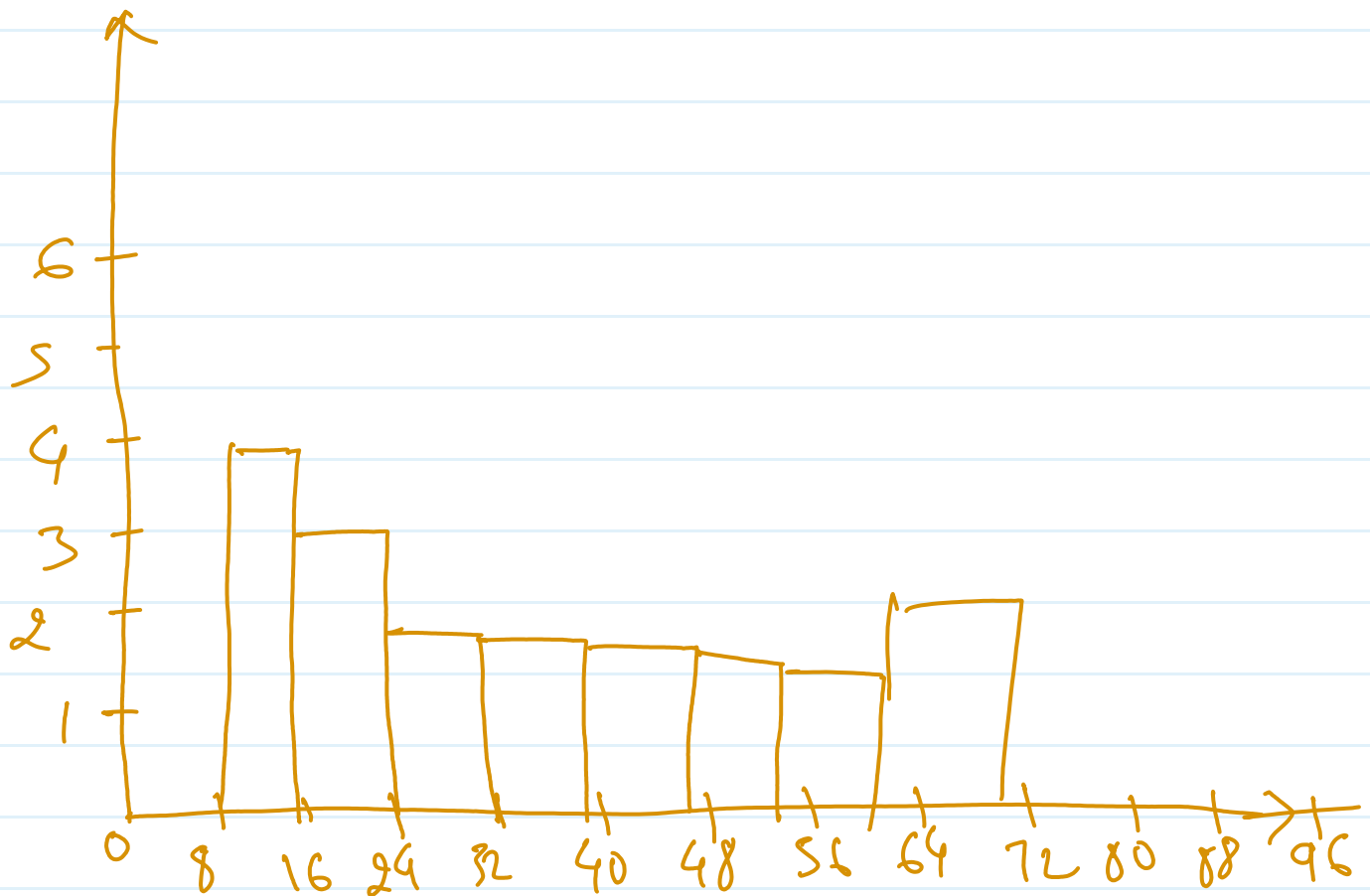
Dataset = [10, 12, 13, 14, 20, 22, 24, 25, 26, 35, 38, 42, 47, 55, 56, 68, 69, 82, 87, 92, ]

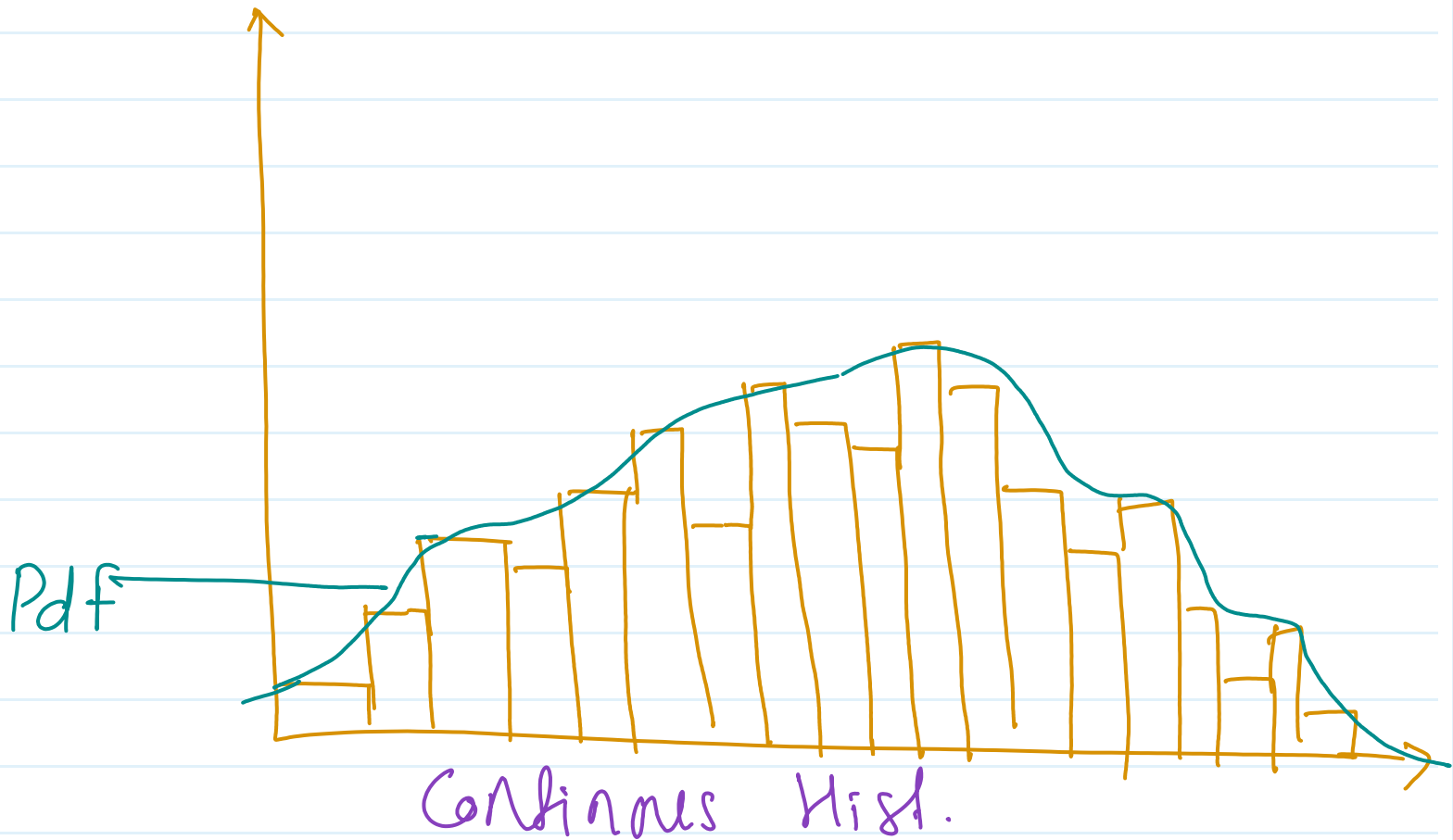
# Bin size

$$\text{Bin size} = \frac{\text{last ele} - \text{first ele.}}{\text{bin no.}}$$

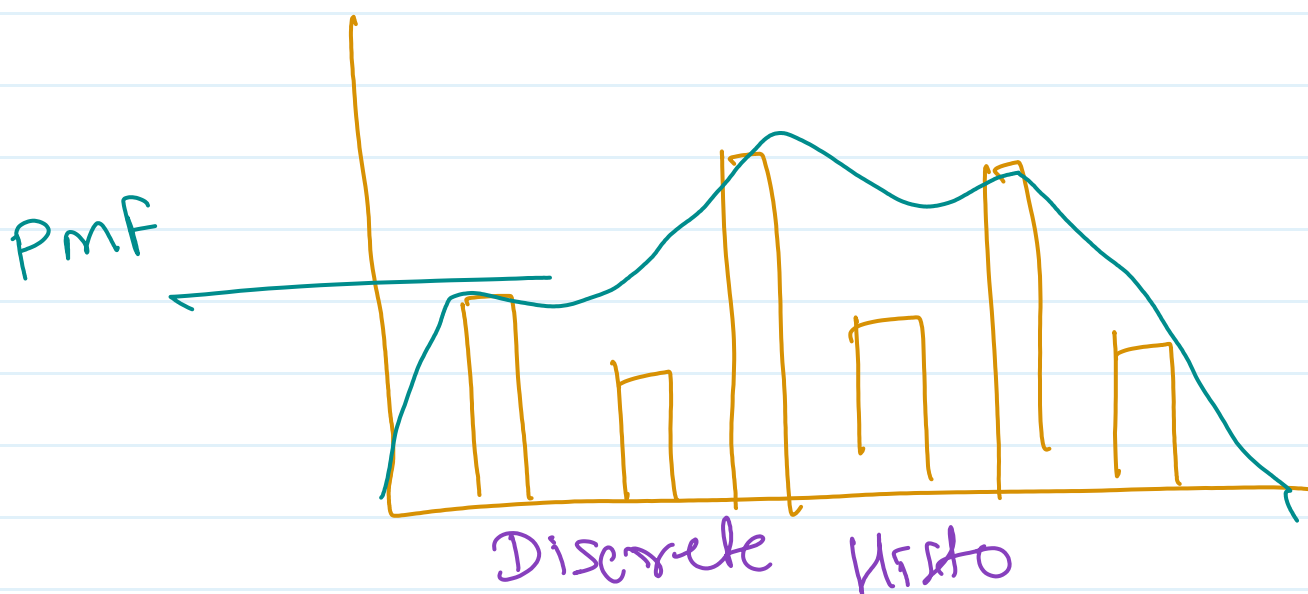
$$= \frac{92 - 10}{10}$$

$$= 82/10 = \underline{8.2}$$





$\text{pdf} = \text{Probability density funct.}$



$\text{Pmf} - \text{Probability mass function.}$



# ★ Covariance

E.g.

X	Y
Age	weight
14	40
15	45
18	51
20	68
25	74

$X \uparrow \quad Y \uparrow$   
 $X \downarrow \quad Y \downarrow$

} +ve covariance

$X \uparrow \quad Y \downarrow$   
 $X \downarrow \quad Y \uparrow$

} -ve covariance

$\begin{matrix} X \uparrow & Y \uparrow \\ X \uparrow & Y \downarrow \\ X \uparrow & Y \uparrow \end{matrix} \left. \vphantom{\begin{matrix} X \uparrow & Y \uparrow \\ X \uparrow & Y \downarrow \\ X \uparrow & Y \uparrow \end{matrix}} \right\} \text{Zero Covariance}$

Age	weight	height	BMI
$X_1$	$X_2$	$X_3$	$Y$
.	.	.	.
—	—	—	9

$X$  = feature column / Independent var.

$Y$  = target column / Dependent var.

Covariance:-

Quality the relationship  
b/w  $X$  &  $Y$  numeric value.

population 
$$\text{Cov}(X, Y) = \sum_{i=1}^N \frac{(X - \bar{X})(Y - \bar{Y})}{N}$$

$$\text{Cov}(X, Y) = \sum_{i=1}^n \frac{(X - \bar{X})(Y - \bar{Y})}{n-1}$$

Eg. Eco. growth                      Altitude

2.1	8
2.5	12
3.6	10
4.0	14

### Covariance

X	Y	$\bar{X}$	$\bar{Y}$	$(X - \bar{X})$	$(Y - \bar{Y})$
2.1	8	3.05	11	-1	-3
2.5	12			-0.6	1
3.6	10			0.5	-1
4.0	14			0.9	3

$$\begin{aligned} \text{Cov}(X, Y) &= \sum_{i=1}^n \frac{(X - \bar{X})(Y - \bar{Y})}{n-1} \\ &= \frac{(-1)(-3) + (-0.6)(1) + (0.5)(-1) + (0.9)(3)}{4-1} \end{aligned}$$



$$= \frac{4.6}{3} = \boxed{\underline{\underline{1.533}}}$$

if value comes in +ve then  
there are positive covariance b/w  
x and y.

X ↑    Y ↑

X ↓    Y ↓

+1000 to -1000

Co-relation ⇒

-1 to 1

Relation

ghn / week

# # Pearson correlation Coefficient

X	Y
2.1	8
2.5	12
3.6	10
4.0	19

Formula

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \cdot \sigma_Y}$$

$$\begin{aligned}\sigma_X &= \sqrt{\sum_{i=1}^n \frac{(X - \bar{X})^2}{n-1}} \\ &= \sqrt{\frac{(-1)^2 + (-0.6)^2 + (0.9)^2 + (0.5)^2}{3}} \\ &= \sqrt{0.8060} = 0.8981\end{aligned}$$

$$\sigma_Y = \underline{\underline{2.58}}$$

$$\text{Cov}(x, y) = \underline{1.533}$$

$$\rho(x, y) = \frac{1.533}{(0.89) \cdot (2.58)}$$

$$= 0.66$$

$$= 66\%$$

0.9 0.4 0.3 0.2 0.1 0 0.1 0.2 0.3 0.9

-0.6  
S

-0.1  
w

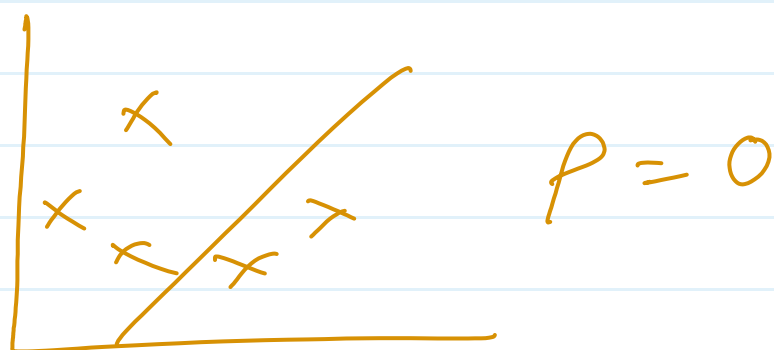
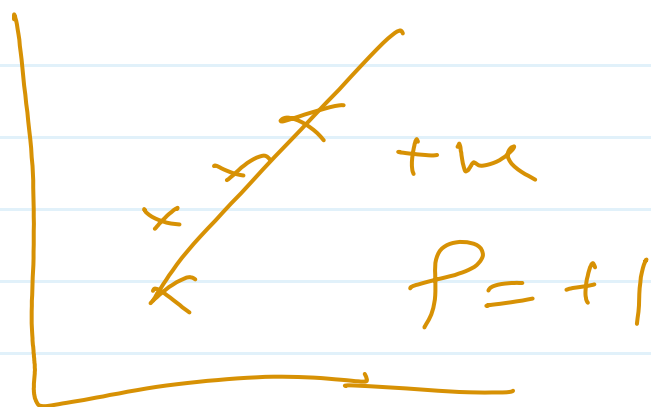
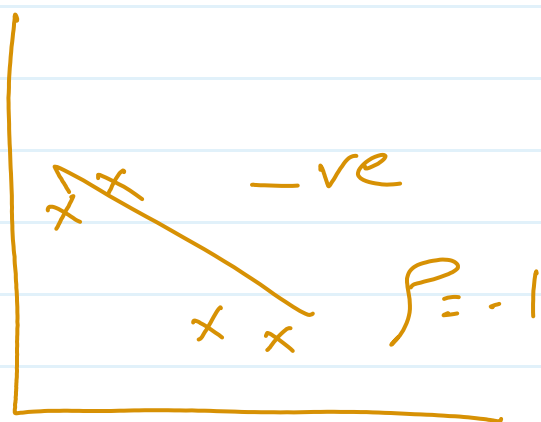
+0.1  
w

+0.7  
S

-0.01

+0.02

✓  
near to zero



### \* Spearman Rank Correlation

$$r_s = \frac{\text{Cov}(R_x, R_y)}{\sigma_{R_x} \sigma_{R_y}}$$

X	Y
2.1	8
2.5	12
3.6	10
4.0	14

Rank  $x$ 

2.1	-	4
2.5	-	3
3.6	-	2
4.0	-	1

Rank  $y$ 

8	-	4
12	-	2
10	-	3
14	-	1

$$\text{mean } R\bar{x} = \frac{4+3+2+1}{4} = 2.5$$

$$\text{mean } R\bar{y} = \frac{4+3+2+1}{4} = 2.5$$

$$(x - \bar{x})$$

$$4 - 2.5 = 1.5$$

$$3 - 2.5 = 0.5$$

$$2 - 2.5 = -0.5$$

$$1 - 2.5 = -1.5$$

$$(y - \bar{y})$$

$$4 - 2.5 = 1.5$$

$$2 - 2.5 = -0.5$$

$$3 - 2.5 = 0.5$$

$$1 - 2.5 = -1.5$$

$$\begin{aligned}
 \text{Cov}_{(x,y)} &= \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1} \\
 &= \frac{(1.5)(1.5) + (0.5)(-0.5) + (-0.5)(0.5) + (-1.5)(-1.5)}{4-1}
 \end{aligned}$$

$$=) \quad 1.33$$

$$\begin{aligned}
 \text{SD } \sigma_x &= \sqrt{\frac{(1.5)^2 + (0.5)^2 + (-0.5)^2 + (-1.5)^2}{4-1}} \\
 &= 1.288
 \end{aligned}$$

$$\text{SD } \sigma_y = 1.288$$

$$\begin{aligned}
 r_s &= \frac{R_x R_y}{\sigma_{R_x} \sigma_{R_y}} \\
 &= \frac{1.33}{1.288 \times 1.288}
 \end{aligned}$$

$$= 0.81$$

$$\boxed{= 81\%}$$

positive correlation is 81% out of  
It is quite strong.