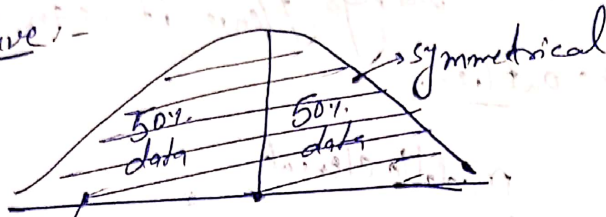


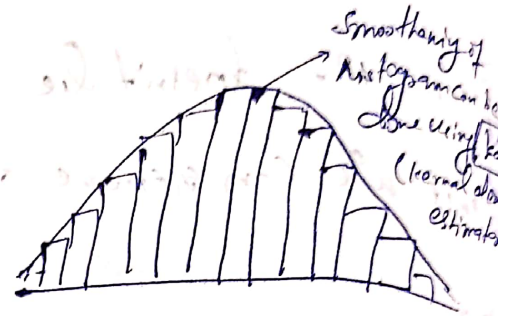
Day 3 Statistics

- ① Normal Distribution / Gaussian Distribution
- ② Standard Normal Distribution
- ③ Z-Score
- ④ Standardization and Normalization
- ⑤ Gaussian / Normal Distribution →

Bell curve :-



Area inside this entire curve is 1 i.e 100%



Most of the data like

Age, weight, height they follow this kind of distribution.

↑
Domain Expertise ⇒ Doctors

IRIS DATASET

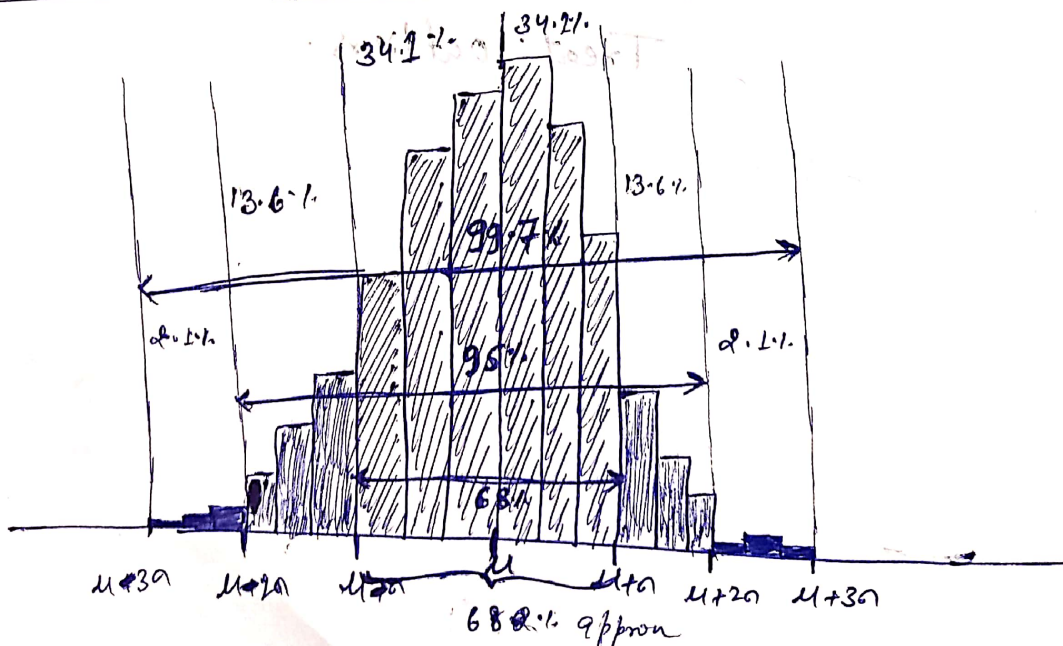


Petal length, sepal length,
Petal width, sepal width

⇓ follow
Gaussian Dist

1.1.2

⇒ Empirical Rule of Normal Distribution :-



Empirical formula (68-95-99.7%) :-

~~68-95-99.7%~~

Empirical formula (68-95-99.7%) rule :-

68% of Entire data present in $(\mu - \sigma, \mu + \sigma)$ region.

95% of Entire data present in $(\mu - 2\sigma, \mu + 2\sigma)$ region.

99.7% of Entire data present in $(\mu - 3\sigma, \mu + 3\sigma)$ region.

and remaining all will be the outside the $(\mu - 3\sigma, \mu + 3\sigma)$ region.
and this rule is called (68-95-99.7% rule) or Empirical formula.

Q-Q Plot \Rightarrow How to determine whether the Distⁿ is Gaussian or not? (Will do it in Practice)

Standard Normal Distribution :-

Random var

$X \approx$ Gaussian Distribution (μ, σ)

\Downarrow Why converting?

convert \Downarrow

$Y \approx$ Standard Normal Distⁿ $(\mu=0, \sigma=1)$ \leftarrow how?

\downarrow simple formula

z-score

\downarrow
sample scaling or
standardization.

Eg:- $X = \{1, 2, 3, 4, 5\}$

$$\mu = 3$$

$$\sigma = 1.41$$

$n=1$

\leftarrow Bcz we are going to apply this formula on each & every value.

$$\text{Z-score} = \frac{X_i - \mu}{\sigma/\sqrt{n}}$$

Standard error

(useful while doing Inferential stats)

if $n \neq 1$ then,

$$\text{Z-score} = \frac{X_i - \mu}{\sigma}$$

there are cases where σ can be different than '1' we will see that later



Now Apply z-score on $x = \{1, 2, 3, 4, 5\}$

$$\mu = 3, \sigma = 1.414$$

$$y = \left\{ -1.414, -0.707, 0, 0.707, 1.414 \right\} = \frac{1-3}{1.414} = -1.414$$

$$\mu = 0, \sigma = 1 = \frac{2-3}{1.414} = -0.707$$

$$= \frac{3-3}{1.414} = 0$$

$$= \frac{4-3}{1.414} = 0.707$$

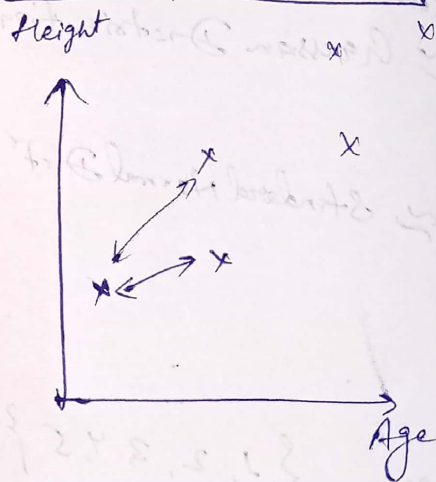
$$= \frac{5-3}{1.414} = 1.414$$

Q. Why are we converting Gaussian Distr to normal standard Distr?

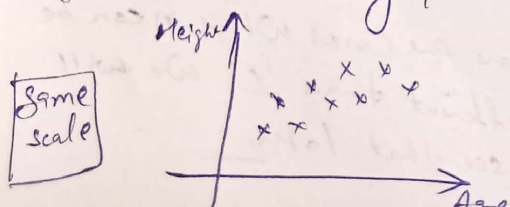
Ans Example \rightarrow

(Years) Age	(kgs) Weight	(cm) Height
24	72	150
26	78	160
32	84	165
33	92	170
34	87	150
28	83	150
29	80	145

Mathematical Calculation Time $\uparrow \uparrow \uparrow$



Hence if we try to bring the data points within the same scale so that it will be nearly populated



Data Points are spread it will take time to calculate the distance between two or for any mathematical calculation it will take time.

So it will have advantage, it will take less time in mathematical calculation. So with this we try to bring data in same scale.

Since by applying z-score formula we bring the data in same scale
i.e. will be ranging b/w -3 to +3 (most of the values) & $\begin{cases} \mu = 0 \\ \sigma = 1 \end{cases}$

But if this calculation will be faster
and we are not losing any data.

And the entire process is called Standardization.

After training the model we can revert the scale back to get the original data. For that again we will have to apply the same formula.

Standardization \rightarrow W/out every ^{feature} value we apply the z-score to scale down the value with $\mu = 0, \sigma = 1$.

$$\text{Age} = \{24, 26, 32, 33, 34, 28, 29\}$$

$$\mu = 29.42$$

$$\sigma^2 = \frac{(24-29.42)^2 + (26-29.42)^2 + (32-29.42)^2 + (33-29.42)^2 + (34-29.42)^2 + (28-29.42)^2 + (29-29.42)^2}{7}$$

$$= \frac{29.37 + 11.69 + 6.65 + 12.81 + 20.97 + 2.016 + 0.176}{7}$$

$$= \frac{83.08}{7}$$

$$= 11.8685$$

$$\sigma = \sqrt{11.8685} = 3.45$$

$$\sigma = 3.45$$

$$\text{Z-Score} = \frac{24-29.42}{3.45}, \frac{26-29.42}{3.45}, \frac{28-29.42}{3.45}, \frac{29-29.42}{3.45},$$

$$\frac{32-29.42}{3.45}, \frac{33-29.42}{3.45}, \frac{34-29.42}{3.45}$$

$$Y_{(\text{Age})} = \{-1.57, -0.99, -0.41, -0.12, 0.74, 1.03, 1.32\}$$

$\mu = 0 \quad \sigma = 1$

Normalization :-

Convert the value b/w any particular range you give that could be $[0, 1]$

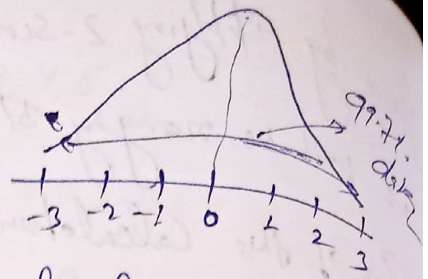
$[0, 1]$ or $[0, 5]$, $[0, 4]$, $[-1, 1]$

Standardization

$$\mu = 0, \sigma = 1$$

$$[-3 \leftrightarrow 3]$$

99.7% of data lies b/w -3 to 3 range



Normalization \rightarrow [Lower Scale \leftrightarrow Higher Scale]

① Min Max Scaler $[0, 1]$

$$x_{\text{Scaled}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Where do we apply this?

in deep learning &

Some of the machine learning techniques

Range \rightarrow Converting \rightarrow apply in CNN

x	y	Calculation
1	0	$\frac{(1-1)}{5-1} = 0$
2	0.25	$\frac{(2-1)}{5-1} = \frac{1}{4} = 0.25$
3	0.5	$\frac{(3-1)}{5-1} = \frac{2}{4} = 0.5$
4	0.75	$\frac{(4-1)}{5-1} = \frac{3}{4} = 0.75$
5	1	$\frac{(5-1)}{5-1} = \frac{4}{4} = 1$

got converted in the range of $[0, 1]$

In Deep Learning :-

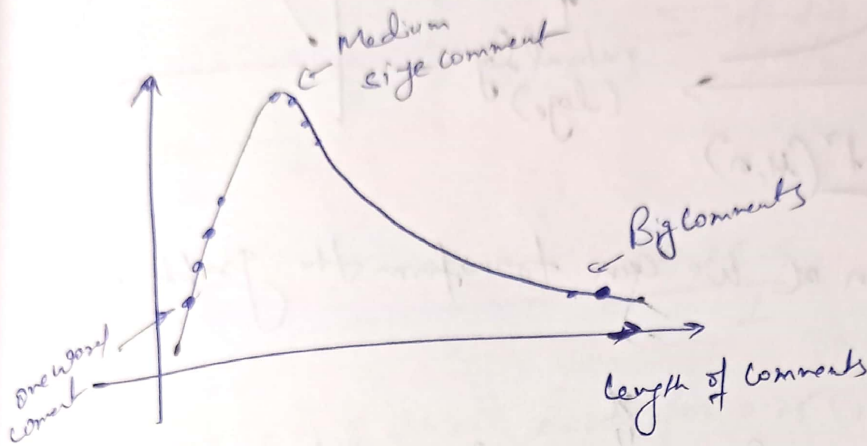
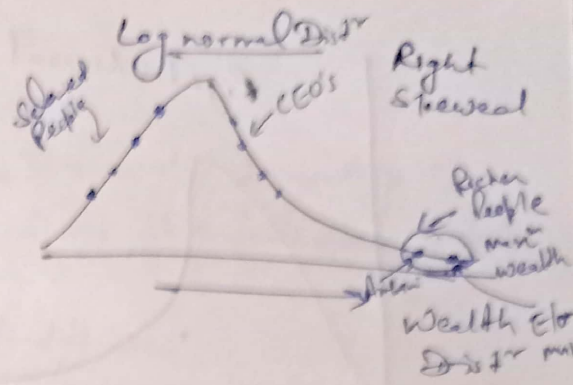
Images \rightarrow Pixels range (0, 255)



\downarrow Convert (scale down)
 $(0, 1) \leftarrow$ Normalization

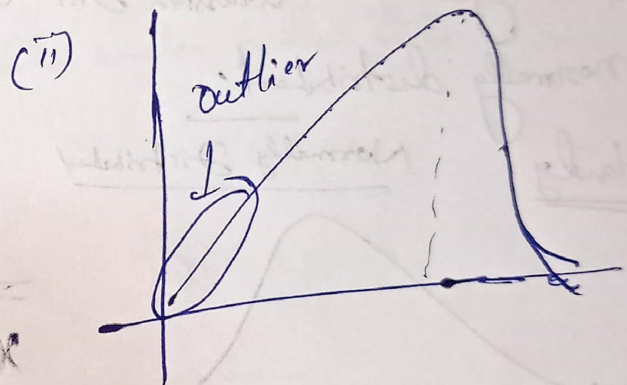
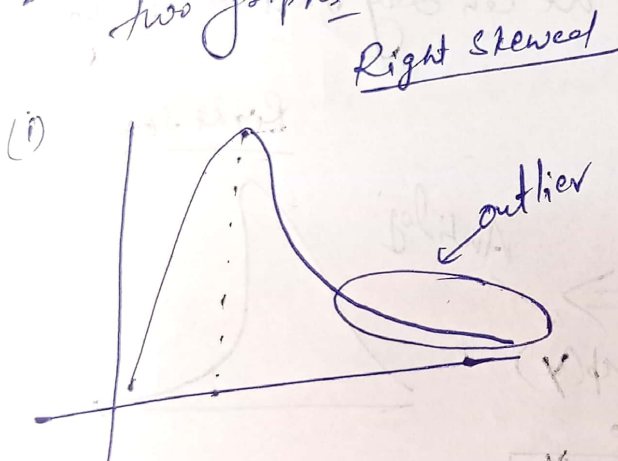
Feature Scaling \leftarrow Standardization
Normalization

① Log Normal Distribution ? →



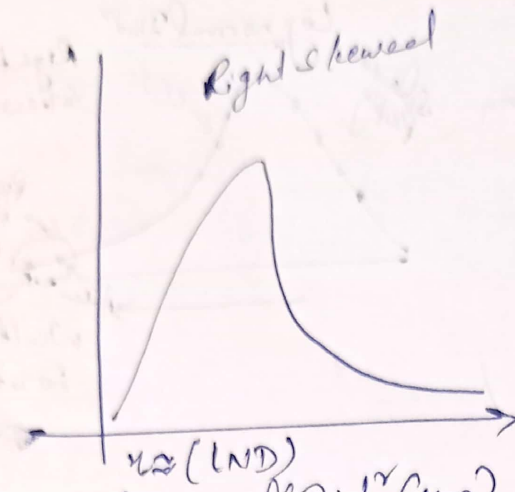
Assignment

Q. What is the relationship of Mean, Median & Mode in this two graphs?



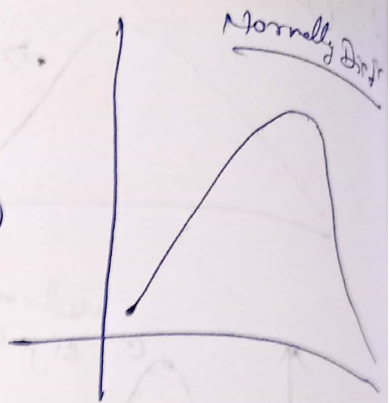
if outlier is bigger
no. it affects the mean
but doesn't affect median.
In (i) outliers are +ve bigger
larger nos. hence mean will be
big. but in case of (ii) outliers are
smaller nos. Hence mean will be small.

~~mean(i) < mean(ii)~~
mean(i) > mean(ii)



$x \sim \text{log-normally Dist}^r(\mu, \sigma)$

$y = \ln(x)$
 \uparrow
 natural log
 (loge)

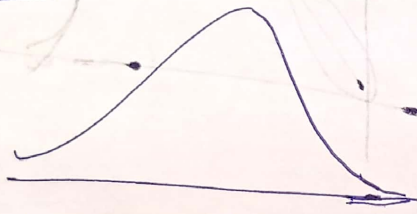


If we perform log on x we can transform the graph in Gaussian Dist^r

\Rightarrow If the random variable x is log normally distributed, then $y = \ln(x)$ has a normal distribution

\Rightarrow If we get y as Gaussian Dist^r then we can say that x is log-normally distributed.

Similarly Normally Distributed

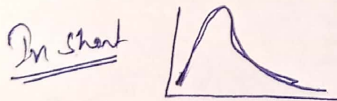


$y \sim \text{normally distributed}$
 $ND(\mu, \sigma)$

\Rightarrow Anti log
 $x = \exp(y)$
 i.e.
 $x = e^y$



has to be as it is normally dist^r

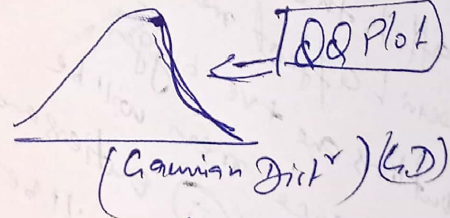


$x \sim \text{LND}$

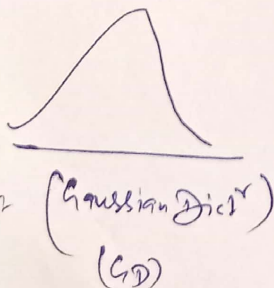
eg. No. of runs

$\Rightarrow y = \log_e(x)$
 or $y = \ln(x)$

\Rightarrow

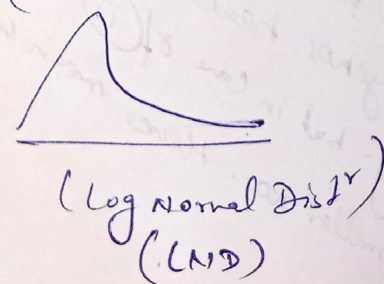


Q-Q Plot

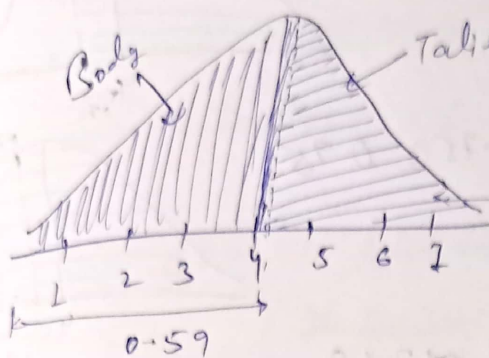


$\Rightarrow y = \exp(x)$
 $y = e^x$

\Rightarrow



Q. $\mu = \{1, 2, 3, 4, 5, 6, 7\}$ $\mu = 4$
 $\sigma = 1$



What is the percentage of Score that falls above 4.25?

find the area under this curve

Sol

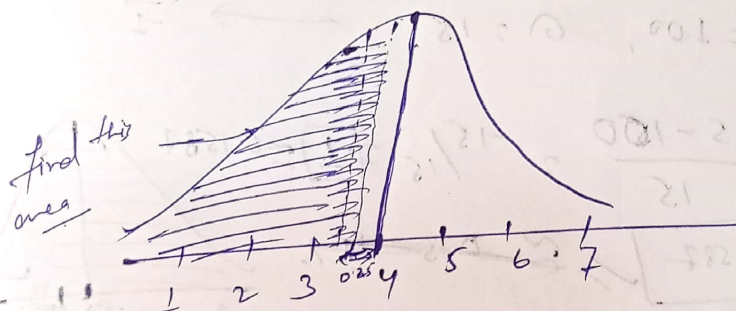
$$Z\text{-Score} = \frac{x_i - \mu}{\sigma} = \frac{4.25 - 4}{1} = \boxed{0.25}$$

now look into Z table look for 0.25 i.e 0.59

Hence area under above 4.25 = ~~0.59~~ $1 - 0.59 = 0.41$

41% Ans

Q. Area under curve that falls under below 3.75?



Sol

$$Z\text{-Score} = \frac{3.75 - 4}{1} = \boxed{-0.25}$$

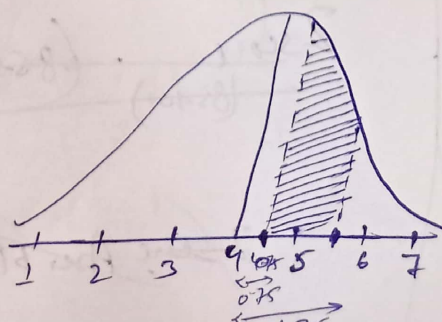
0.25 away from the mean

now look into Z-table & look for -ve table for -0.25 i.e 0.40
40% Ans

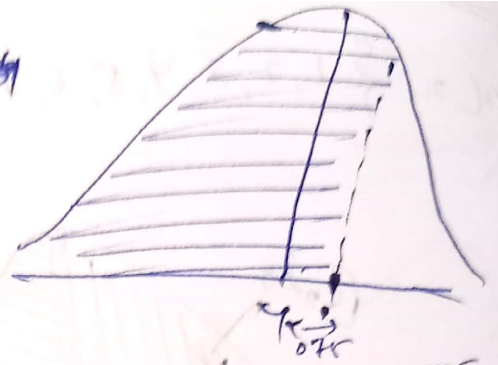
Q. Find the area of curve b/w 4.75 & 5.75

$$Z\text{-Score} = \frac{4.75 - 4}{1} = 0.75$$

$$Z\text{-Score} = \frac{5.75 - 4}{1} = 1.75$$



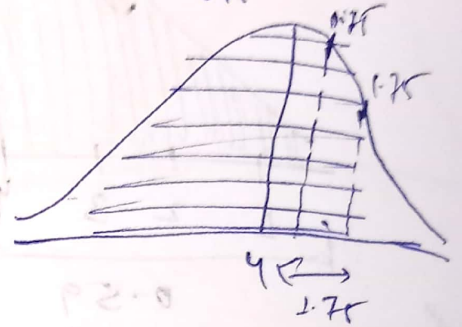
Now look into Z table for $0.75 \approx 0.7734$



Now look into Z table for $1.75 \approx 0.9599$

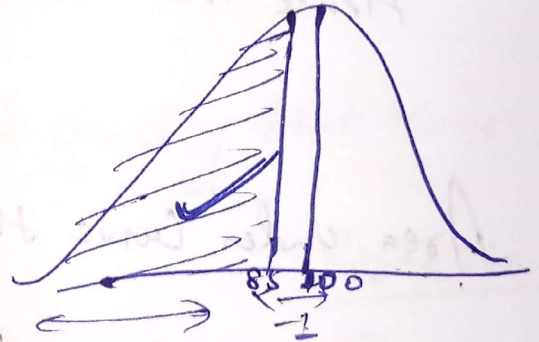
Now subtract $Z_{1.75} - Z_{0.75}$

$$(0.9599 - 0.7734) \approx 0.1865 = 18.65\% \text{ Ans}$$



Q In India the Avg IQ is 100 with a S.D of 15 what is the percentage of population could you expect to have an IQ

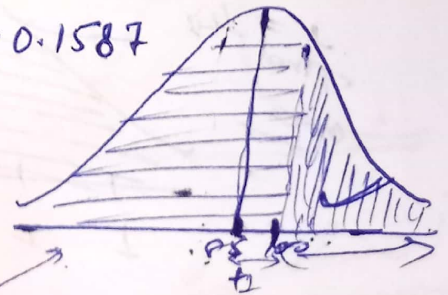
- lower than 85 = 0.1587
- higher than 85 = 0.8413
- Between 85 & 100 = 0.3413



Sol
 $\mu = 100, \sigma = 15$

$$Z\text{score} = \frac{85 - 100}{15} = \frac{-15}{15} = -1 \approx 0.1587$$

Area below 85 ≈ 0.1587



Higher than 85 = $1 - 0.1587 = 0.8413$

Area above 85

Between 85 & 100

Area ^{higher} than 85 = ~~0.8413~~
_{lower} 0.1587

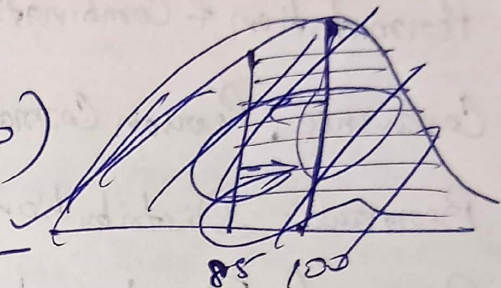
$z \times 100 = 0 \approx \cancel{0.5000}$



~~Area~~ Area below 100 = 0.5000

Now Area b/w 85 & 100 (~~0.8413 - 0.5000~~)

~~0.3413~~



(Area below 100 - Area below 85)

$(0.5000 - 0.1587) \approx 0.3413$

$85 < \text{area} < 100$

0.3413

