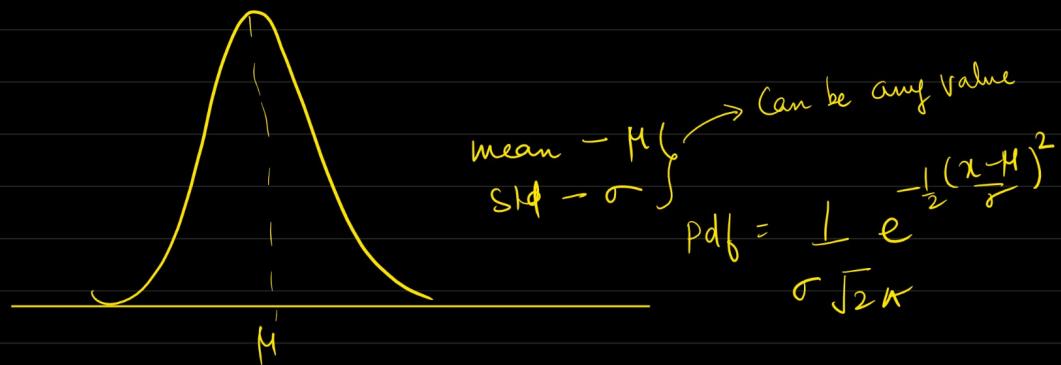
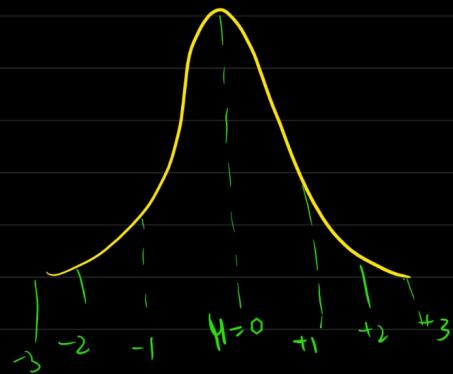


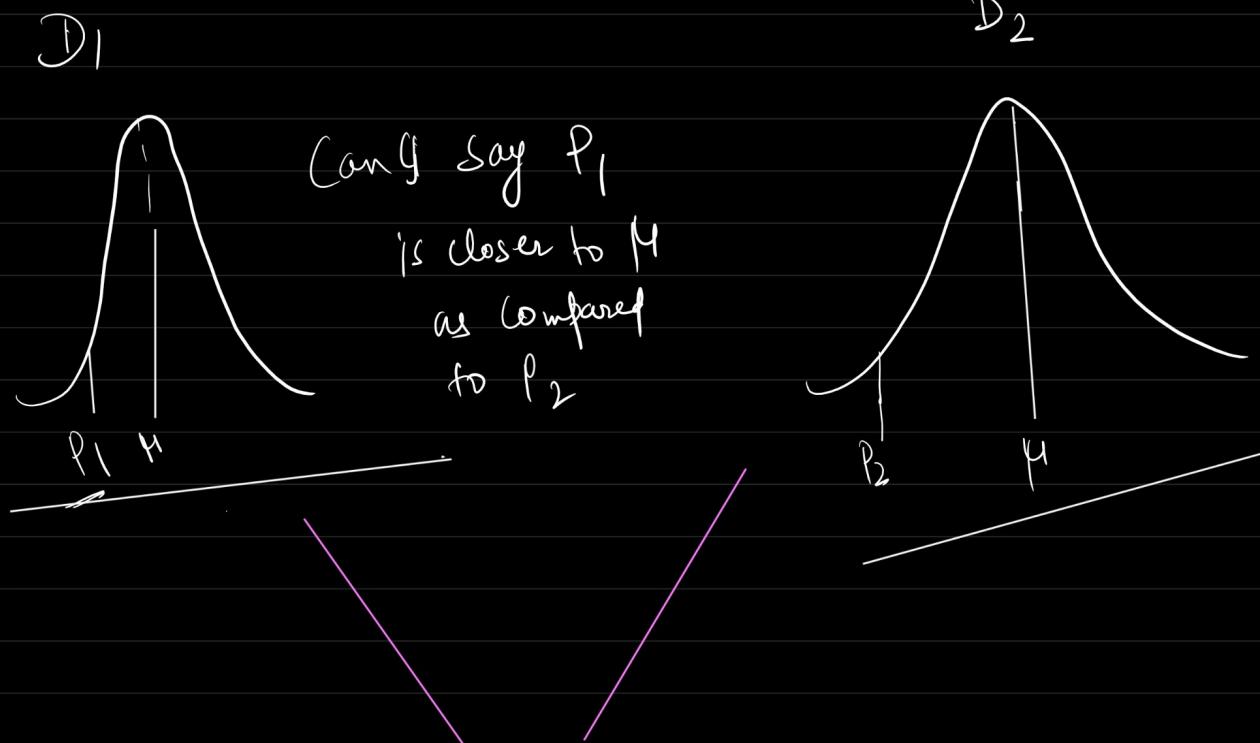
Standard Normal distribution



- S.N.D. is continuous prob distribution
- A special case of N.D. $\rightarrow \mu=0, \sigma=1$



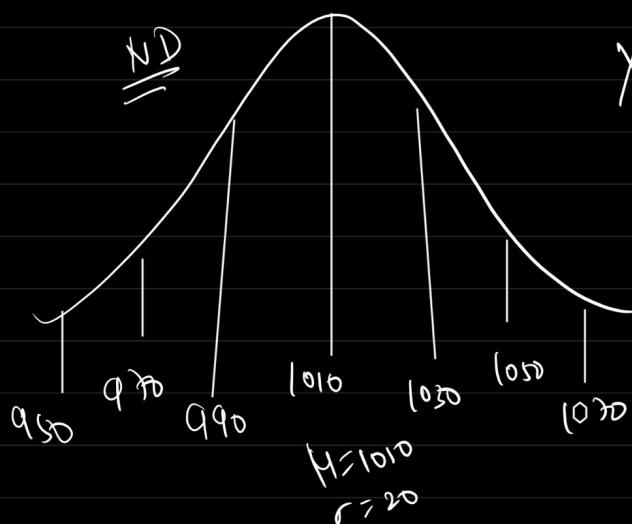
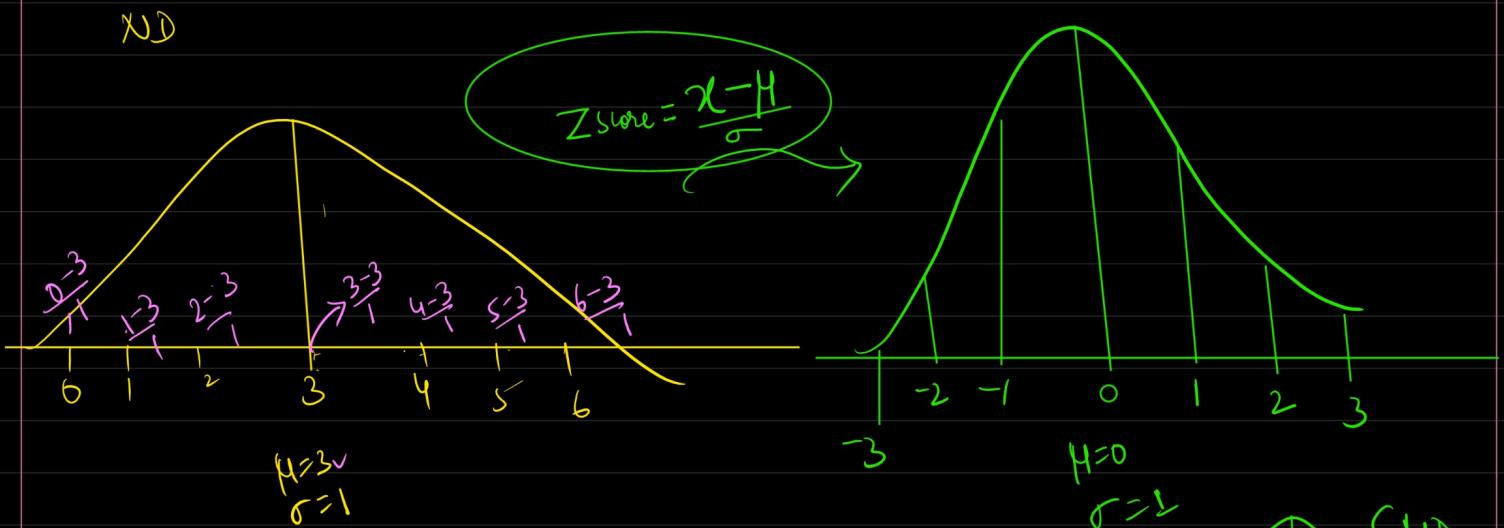
why SND if we already have N.D?



On the same scale $\bar{M} \approx \sigma$

$N \cdot D$

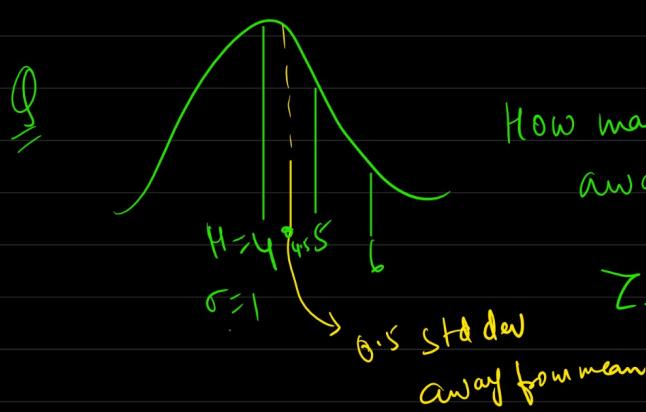
Standard $N \cdot D$



$$X \approx SND (M=0, \sigma=1)$$

$$Z = \frac{x-M}{\sigma}$$

x	$x-M$	$\frac{x-M}{\sigma}$
950	$950 - 1010 = 60$	$\frac{60}{20} = 3$
970	$970 - 1010 = 40$	2
990	$990 - 1010 = 20$	1
1010	$1010 - 1010 = 0$	0
1030	$1030 - 1010 = 20$	1
1050	$1050 - 1010 = 40$	2
1070	$1070 - 1010 = 60$	3



How many standard deviation 4.5 is away from mean?

$$Z\text{ score} = \frac{x-M}{\sigma} = \frac{4.5 - 4}{1} \Rightarrow 0.5$$

Σ	No of rooms	Area of house (sq ft)	Locality	Distance from airport	Y (Price of house) (in lakhs)
1	1100	1		20 km	70
2	1200	2		30 km	80
3	1150	3		15 km	90
4	1250	1		60 km	110
5	1300	1		30 km	120
-	-	2		-	-
-	-	-		-	-
-	-	-		-	-

Linear Reg
Logistic Reg
Clustering

→ Standardization

$$\rightarrow Z_{\text{Score}} = \frac{x - \bar{x}}{\sigma}$$

$$= \frac{x_i - \text{Area of house}}{\sigma}$$

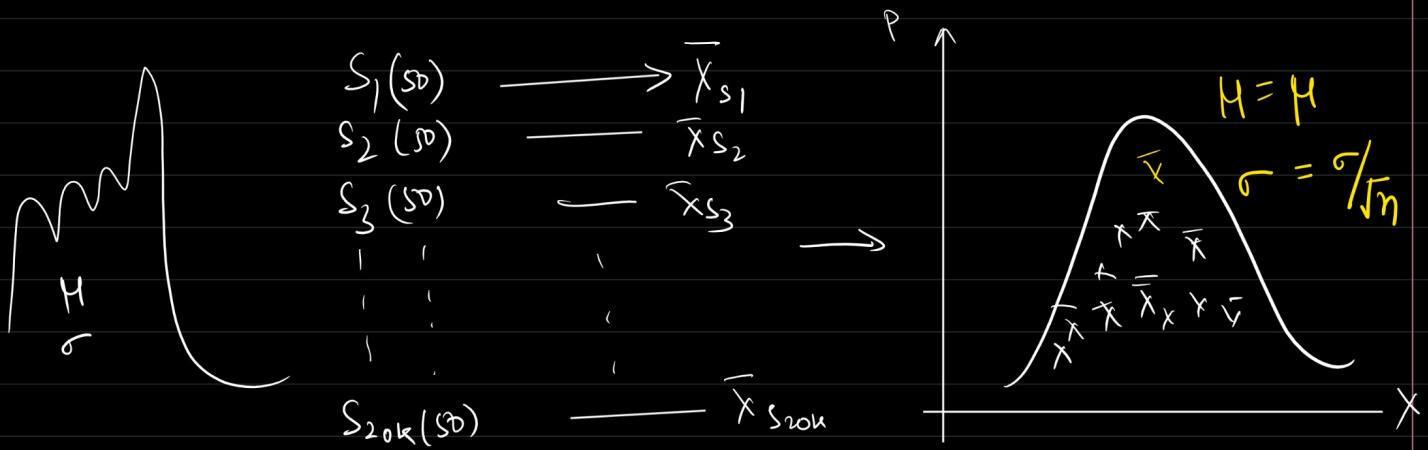
Central limit theorem

→ Distribution will be irregular.



→ The CLT states that if you have a population with mean μ and standard deviation σ and take sufficiently large random samples from the population with replacement, then the distribution of sample means will be approximately normally distribution.

→ Sampling mean of a population (μ, σ) will approximately be a normal distribution $\rightarrow \mu = \mu, \sigma = \sigma/\sqrt{n}$



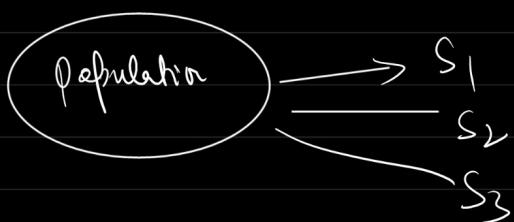
CLT → Sampling distribution of the mean \Rightarrow

Population $(\mu, \sigma) \rightarrow$ Large No of Sample \rightarrow Sampling mean \rightarrow Plot $\hookrightarrow N(1, \mu = \mu, \sigma = \sigma/\sqrt{n})$

Two conditions of CLT :-

- (1) The no. of samples should be large.
- (2) The sample size should be greater than or equals to 30 (except the population distribution which is already a Normal dist.)

$$\text{Standard error} = \sigma / \sqrt{n}$$



$$\frac{\sigma}{\sqrt{n}}$$

Higher the sample size, SE will be low

$$SE \propto \frac{1}{\sqrt{n}} \quad SE \downarrow \quad n \uparrow$$

You have a population with a $\mu = 100$ and Std dev (σ) = 20. If you take sample size 5 from this population. What is the prob that Sample mean will be less than 105.

$$\rightarrow \mu = 100, \sigma = 20, n = 50, \bar{X} = 105$$

$$Z_{\text{score}} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{105 - 100}{20 / \sqrt{50}} = \frac{5 \cdot \sqrt{2}}{4}$$

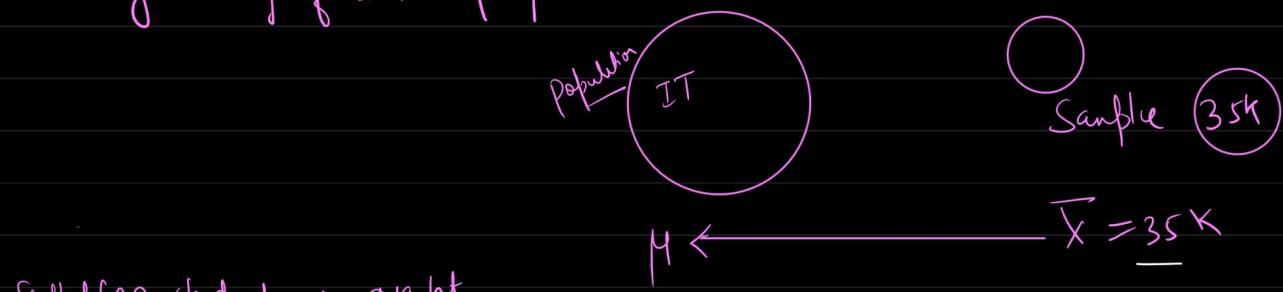
$$Z_{\text{score}} = \frac{5 \cdot \sqrt{2}}{4} \Rightarrow \frac{5 \cdot 1.414}{4} = \frac{7.07}{4} = 1.7675 \\ = 0.8944$$



Estimate → A specified observed value used to estimate an unknown population parameter using sample.

① Point estimate → A point estimate is a single value that is used to estimate the true value of a population parameter.

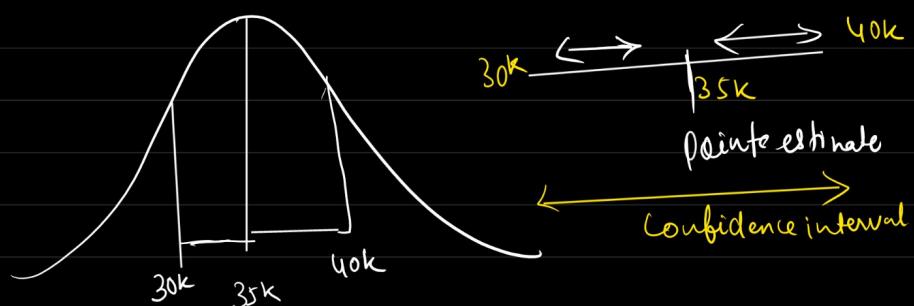
Arg salary of IT employees



Sample of 100 students → avg ht

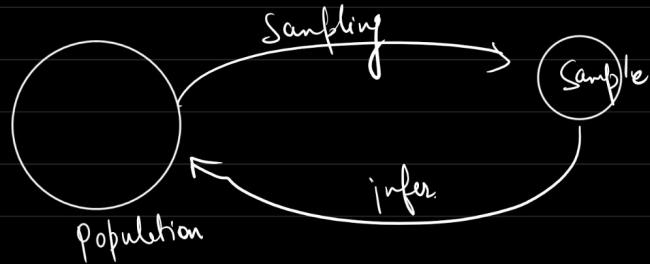
→ mean, variance

② Interval estimate → Range of values used to estimate the unknown population parameters



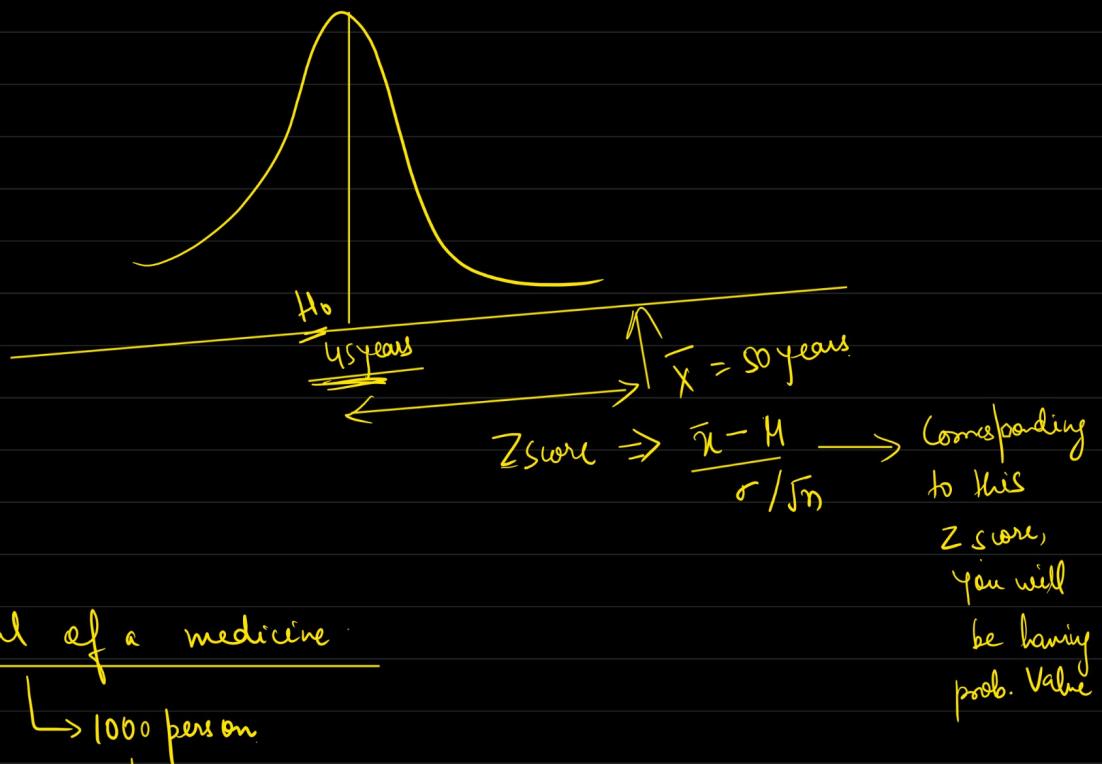
P-value

- The p-value is the probability value, calculated from a statistical test.
- P-value in Hypothesis testing is used to decide whether to reject a null hypothesis or not.



* Age of the employee is 45 years (pop)

- You took a sample
- Calculated avg age of that sample → 50 years.



Clinical trial of a medicine

↳ 1000 person

↳ 95 → got cured with this medicine.

↳ Medicine is working 95%.

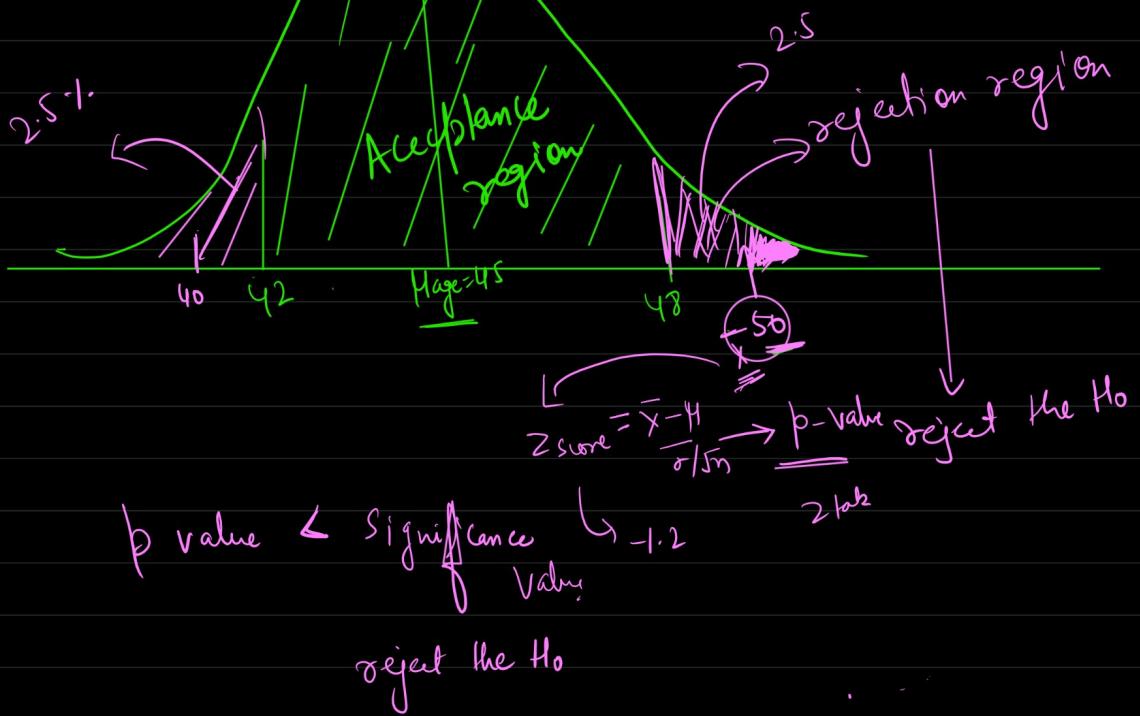
→ Out of 100 times, 95% will work. → [95% Confident]

→ 5% of the time medicine doesn't work. — [5% margin of error]

* Conduct experiment with 5% level of significance

$$CI = 1 - 0.05 \\ \Rightarrow 0.95$$

5%



else fail to reject the H_0

Application of Z score

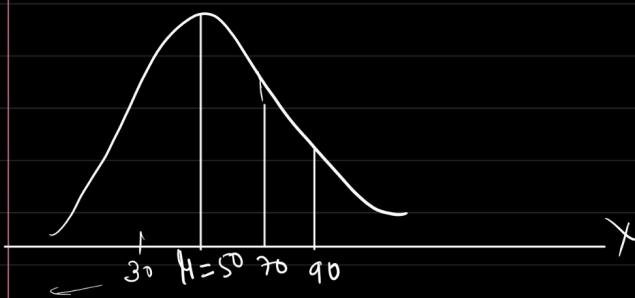
$$Z \text{ score} = \frac{x - \mu}{\sigma}$$

$$\therefore N(\mu = 50 \text{ cm}, \sigma = 20, D = 110)$$

How many standard deviation D is away from mean?

$$Z \text{ score} = \frac{x_i - \mu}{\sigma} = \frac{110 - 50}{20} = \frac{60}{20} = 3$$

$D = 110$ is 3σ away from mean.



$Z = 3 \rightarrow D = 110$ is 3 SD away from mean
 $\sigma = 20 \rightarrow$ the dp on an avg is 20 units away from mean

$$1 \text{ SD away from mean} \Rightarrow \mu + 1\sigma = 50 + 1 \times 20 = 70$$

$$2 \text{ " " } \Rightarrow \mu + 2\sigma = 50 + 2 \times 20 \Rightarrow 90$$

$$3 \text{ SD " " " } \Rightarrow \mu + 3\sigma = 50 + 3 \times 20 \Rightarrow 110$$

$$\begin{aligned} &= \mu - 1\sigma \\ &= \mu - 2\sigma \\ &= \mu - 3\sigma \end{aligned}$$

$$\begin{aligned} \mu - 3\sigma &= 4 - 3 \times 1 = 1 \\ \mu - 2\sigma &= 4 - 2 \times 1 = 2 \\ \mu - 1\sigma &= 4 - 1 = 3 \end{aligned}$$

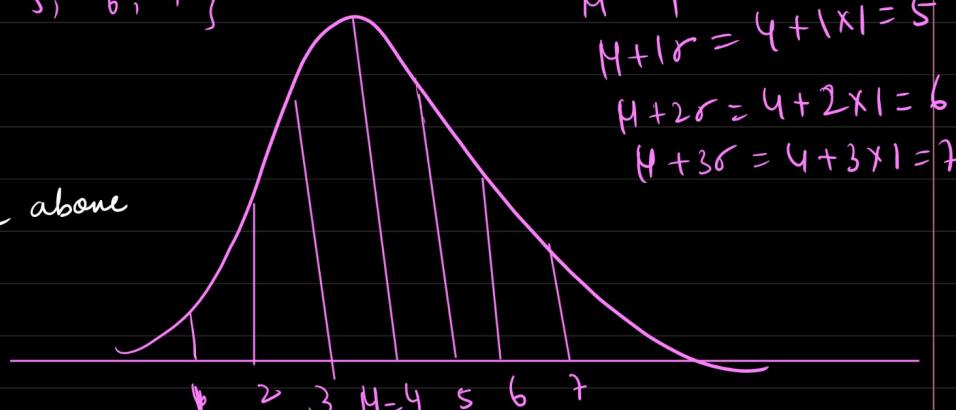
$$\therefore X = \{1, 2, 3, 4, 5, 6, 7\}$$

$$\mu = 4$$

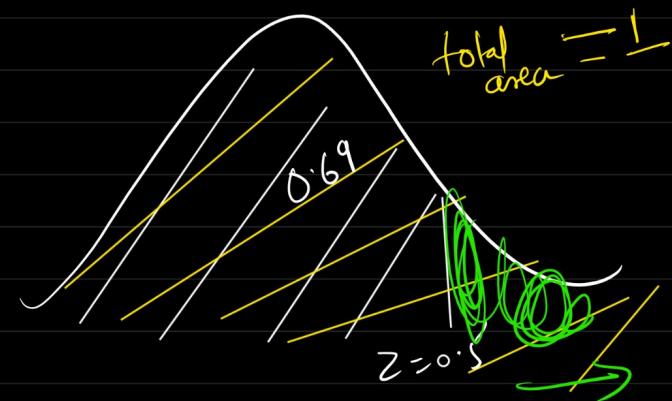
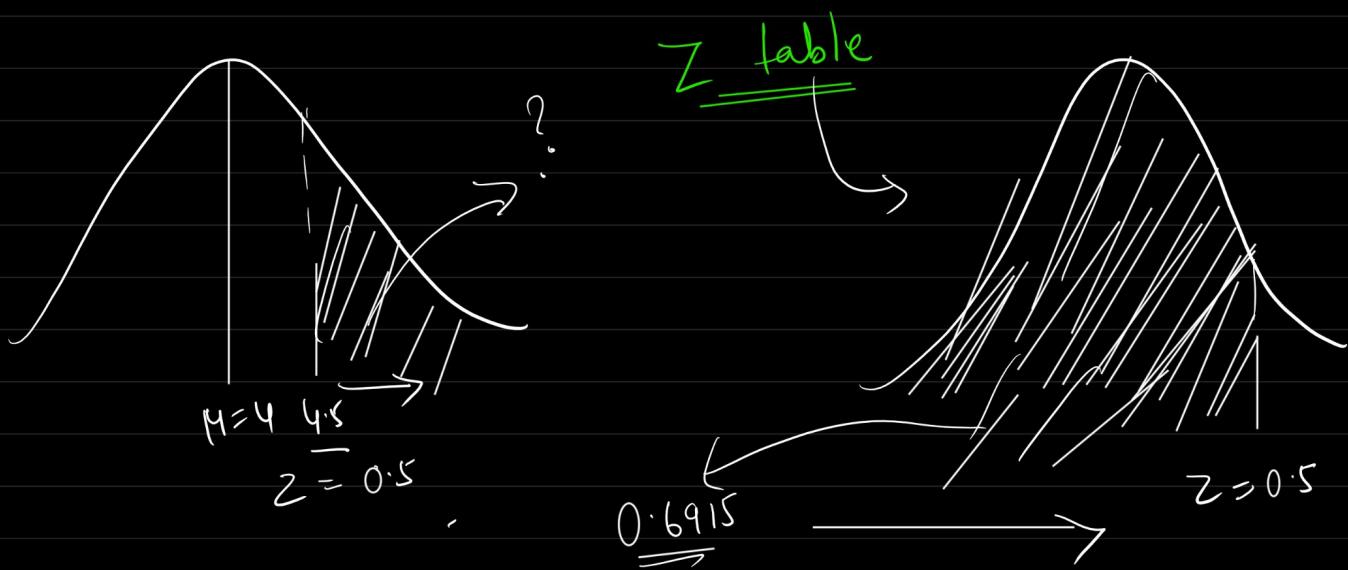
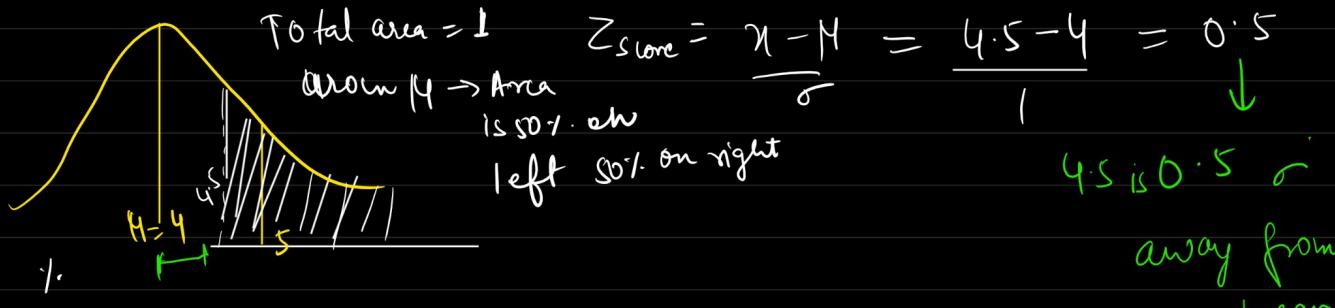
$$\sigma = 1$$

What % of score will fall above 4.5?

→ What is probability that score is more than 4.5?



$$\begin{aligned} \mu &= 4 \\ \mu + 1\sigma &= 4 + 1 \times 1 = 5 \\ \mu + 2\sigma &= 4 + 2 \times 1 = 6 \\ \mu + 3\sigma &= 4 + 3 \times 1 = 7 \end{aligned}$$



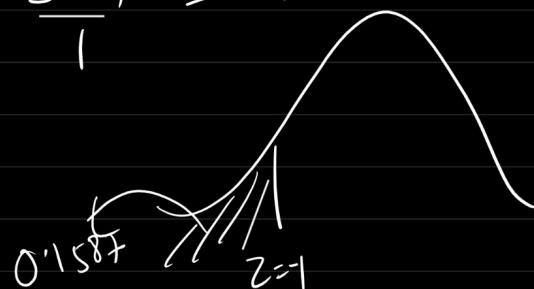
$$= 1 - 0.69 \\ = 0.31$$

Q What is Percentage of marks below 3?

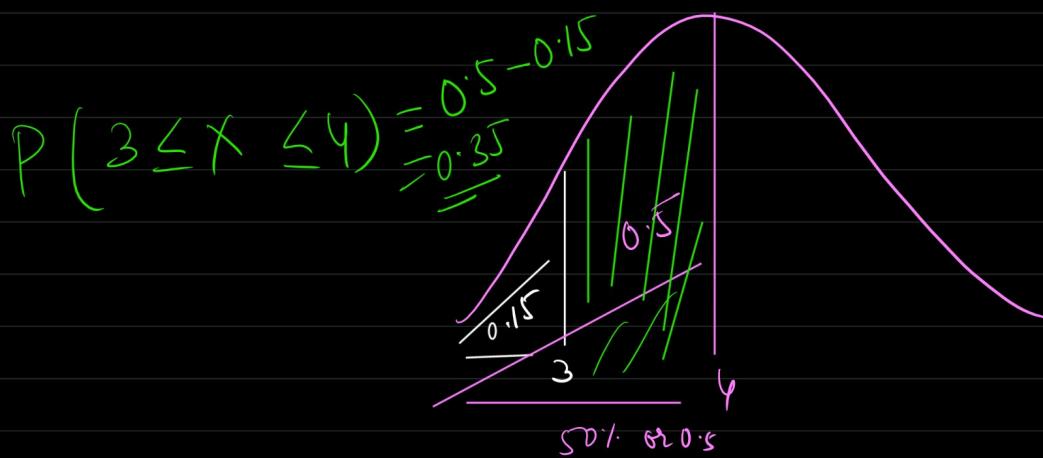
$$Z = \frac{x-\mu}{\sigma} = \frac{3-4}{1} = -1$$

$$\underline{Z = -1}$$

$$0.1587$$



Q What is the age of marks between 4 & 3



Q The score follows a SND $\mu = 75$, and $\sigma = 10$

find prob that a randomly selected student will score below 80.

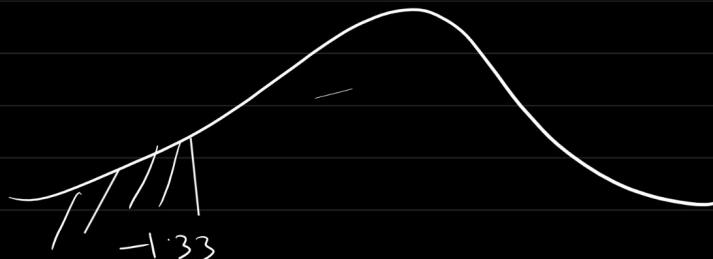
$$Z = \frac{x - \mu}{\sigma} = \frac{80 - 75}{10} = \frac{1}{2} \Rightarrow 0.5$$

for $Z = 0.5$, AVG/cumulative prob = 0.6915

Q The avg IQ is 100 with $\sigma = 15$

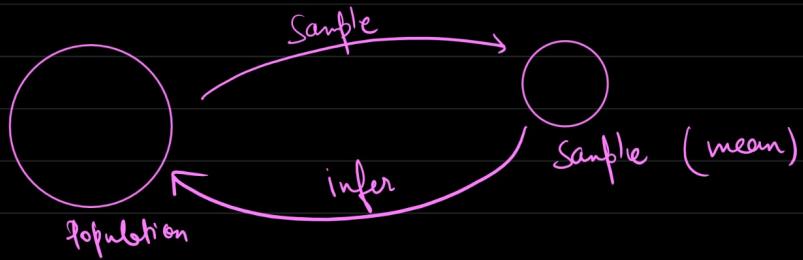
What is the prob of people lower than IQ 80

$$\rightarrow Z_{\text{score}} = \frac{80 - 100}{15} = \frac{-20}{15} = -1.33$$



$$\Rightarrow 0.9082$$

Hypothesis and mechanism



→ Hypothesis testing.

* Hypothesis → It is a claim or a statement or an assumption about a population parameter that can be tested using statistical methods.

e.g. Avg Salary of IT employee \$0K.

e.g. Consumption of Ice-Cream is more in Summer.

e.g. The person is not guilty if accused of any crime.

① Null hypothesis → The initial or default assumption.

e.g. Person is not guilty

② Alternate Hypothesis → Opposite of Null Hypothesis

e.g. Person is guilty.

Mechanism

① Frame the hypothesis

↳ H_0 (Null hypothesis)
↳ H_A (Alternate hypothesis)

Claim hypothesis :- Avg age of people in PwSkills is 45 years.

$$H_0 : \text{Avg age} = 45$$

$$H_A : \text{Avg age} \neq 45$$

$H_0 \rightarrow$ will have equality sign.

Hypothesis Avg age of Employee in ABC organisation is atleast 45 years?

$$\begin{cases} H_0: \text{Hage} \geq 45 \\ H_A: \text{Hage} < 45 \end{cases}$$

$$\begin{cases} \text{Hage} \geq 45 \\ \text{Hage} < 45 \end{cases}$$

→ Avg age of Employee in pwskills is greater than 45 years?

$$\begin{cases} \text{Hage} > 45 \\ \text{Hage} \leq 45 \end{cases}$$

$$H_0: \text{Hage} \leq 45$$

$$H_A: \text{Hage} > 45$$

→ Avg age is greater than equals 45 years

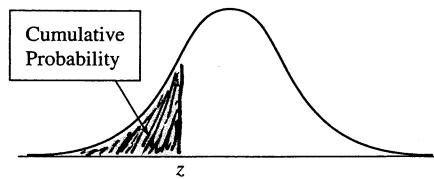
$$H_0: \text{Hage} \geq 45$$

$$H_A: \text{Hage} < 45$$

② Experiment / statistical analysis (pvalue, significance level)

③ reject H_0 or fail to reject H_0

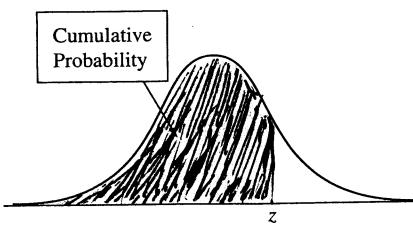
APPENDIX A



Cumulative probability for z is the area under the standard normal curve to the left of z

TABLE A Standard Normal Cumulative Probabilities

z	.00	z	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
-5.0	.000000287	-3.4	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0002
-4.5	.00000340	-3.3	.0005	.0005	.0005	.0004	.0004	.0004	.0004	.0004	.0004	.0003
-4.0	.0000317	-3.2	.0007	.0007	.0006	.0006	.0006	.0006	.0006	.0005	.0005	.0005
<u>-3.5</u>	<u>.000233</u>	-3.1	.0010	.0009	.0009	.0009	.0008	.0008	.0008	.0008	.0007	.0007
		-3.0	.0013	.0013	.0013	.0012	.0012	.0011	.0011	.0011	.0010	.0010
		-2.9	.0019	.0018	.0018	.0017	.0016	.0016	.0015	.0015	.0014	.0014
		-2.8	.0026	.0025	.0024	.0023	.0023	.0022	.0021	.0021	.0020	.0019
		-2.7	.0035	.0034	.0033	.0032	.0031	.0030	.0029	.0028	.0027	.0026
		-2.6	.0047	.0045	.0044	.0043	.0041	.0040	.0039	.0038	.0037	.0036
		-2.5	.0062	.0060	.0059	.0057	.0055	.0054	.0052	.0051	.0049	.0048
		-2.4	.0082	.0080	.0078	.0075	.0073	.0071	.0069	.0068	.0066	.0064
		-2.3	.0107	.0104	.0102	.0099	.0096	.0094	.0091	.0089	.0087	.0084
		-2.2	.0139	.0136	.0132	.0129	.0125	.0122	.0119	.0116	.0113	.0110
		-2.1	.0179	.0174	.0170	.0166	.0162	.0158	.0154	.0150	.0146	.0143
		-2.0	.0228	.0222	.0217	.0212	.0207	.0202	.0197	.0192	.0188	.0183
		-1.9	.0287	.0281	.0274	.0268	.0262	.0256	.0250	.0244	.0239	.0233
		-1.8	.0359	.0351	.0344	.0336	.0329	.0322	.0314	.0307	.0301	.0294
		-1.7	.0446	.0436	.0427	.0418	.0409	.0401	.0392	.0384	.0375	.0367
		-1.6	.0548	.0537	.0526	.0516	.0505	.0495	.0485	.0475	.0465	.0455
		-1.5	.0668	.0655	.0643	.0630	.0618	.0606	.0594	.0582	.0571	.0559
		-1.4	.0808	.0793	.0778	.0764	.0749	.0735	.0721	.0708	.0694	.0681
		-1.3	.0968	.0951	.0934	.0918	.0901	.0885	.0869	.0853	.0838	.0823
		-1.2	.1151	.1131	.1112	.1093	.1075	.1056	.1038	.1020	.1003	.0985
		-1.1	.1357	.1335	.1314	.1292	.1271	.1251	.1230	.1210	.1190	.1170
		-1.0	.1587	.1562	.1539	.1515	.1492	.1469	.1446	.1423	.1401	.1379
		-0.9	.1841	.1814	.1788	.1762	.1736	.1711	.1685	.1660	.1635	.1611
		-0.8	.2119	.2090	.2061	.2033	.2005	.1977	.1949	.1922	.1894	.1867
		-0.7	.2420	.2389	.2358	.2327	.2296	.2266	.2236	.2206	.2177	.2148
		-0.6	.2743	.2709	.2676	.2643	.2611	.2578	.2546	.2514	.2483	.2451
		-0.5	.3085	.3050	.3015	.2981	.2946	.2912	.2877	.2843	.2810	.2776
		-0.4	.3446	.3409	.3372	.3336	.3300	.3264	.3228	.3192	.3156	.3121
		-0.3	.3821	.3783	.3745	.3707	.3669	.3632	.3594	.3557	.3520	.3483
		-0.2	.4207	.4168	.4129	.4090	.4052	.4013	.3974	.3936	.3897	.3859
		-0.1	.4602	.4562	.4522	.4483	.4443	.4404	.4364	.4325	.4286	.4247
		-0.0	.5000	.4960	.4920	.4880	.4840	.4801	.4761	.4721	.4681	.4641



Cumulative probability for z is the area under the standard normal curve to the left of z

TABLE A Standard Normal Cumulative Probabilities (continued)