

## Normal Approximation & Binomial Distribution

→ Empirical rule and normal approximation of data

→ Many data have histograms that look bell-shaped. If data has a normal bell, then it is a normal data normal curve

- ① has bell curve
- ② is symmetrical

e.g. - height, weight, B.P

### The Normal Approximation:-

If data has a normal curve, then

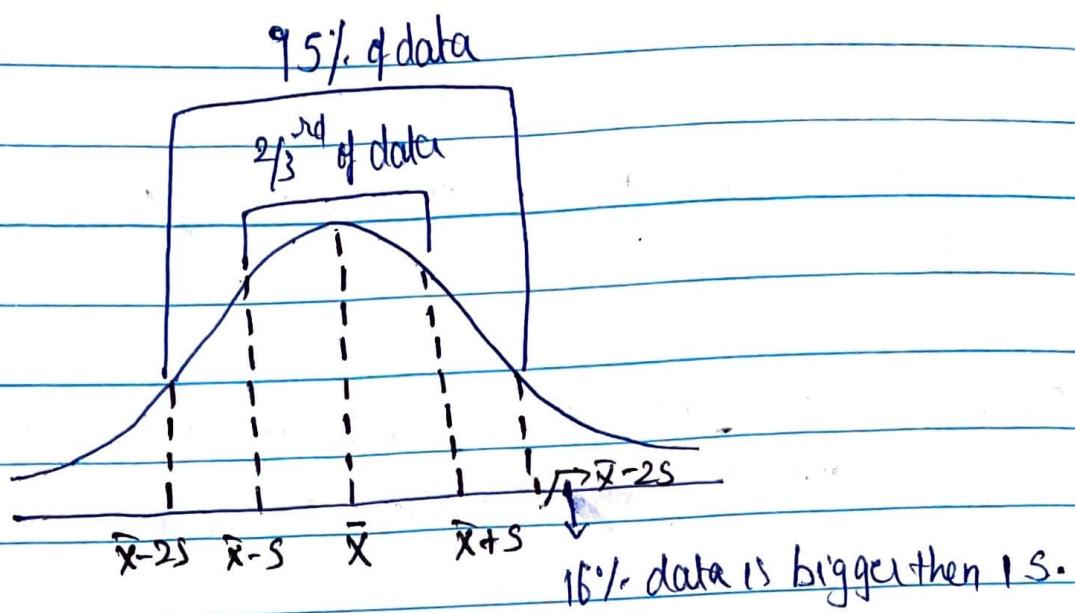
- about  $\frac{2}{3}$  rd of data fall within one standard deviation of mean
- about 95% fall within 2 standard deviations of mean.
- about 99.7% fall within 3 standard deviations.

e.g. if  $s = 1.8$  &  $\bar{x} = 68.3$ , then

$$\bar{x} - 2(1.8) = 64.7$$

$$\bar{x} + 2(s) = 71.9$$

∴ 95% of data lies b/w 64.7 and 71.9.



## Standardizing data & standard normal curve

Normal curve is completely determined by  $\bar{x}$  and  $s$ . Once we know them, we know whole histogram & we can compute  $\%$ . To do that, we need to standardize data.

Standardizing data that we take data by subtracting off  $\bar{x}$  & dividing by  $s$ .

$$Z = \frac{\text{data} - \bar{x}}{s}$$

$Z = Z\text{-score} = \text{standardized value}$ .

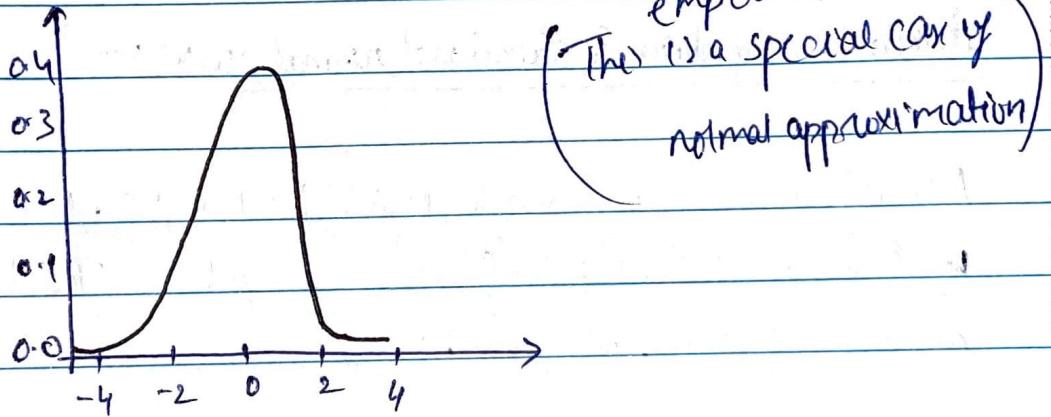
$\Rightarrow$  Z-score tells how many standard deviations, the measurement is above or below average.

e.g:- if  $Z=2$ , i.e data is 2 standard deviations above average.  
if  $Z=-1.5$ , i.e data is 1.5 standard deviations below average.

Once we standardize data that have  $\bar{x}=0$  &  $s=1$ , this is point of standardizing.

If  $\bar{x}=68.3$  &  $s=1.8$ , after we standardize, the standardized values follow "standard normal curve" which means has  $\bar{x}=0$  &  $s=1$ .

Standard normal curve graph,  $y = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}$



### Normal Approximation :-

Normal approximation means to use the area under the normal curve to figure percentages.

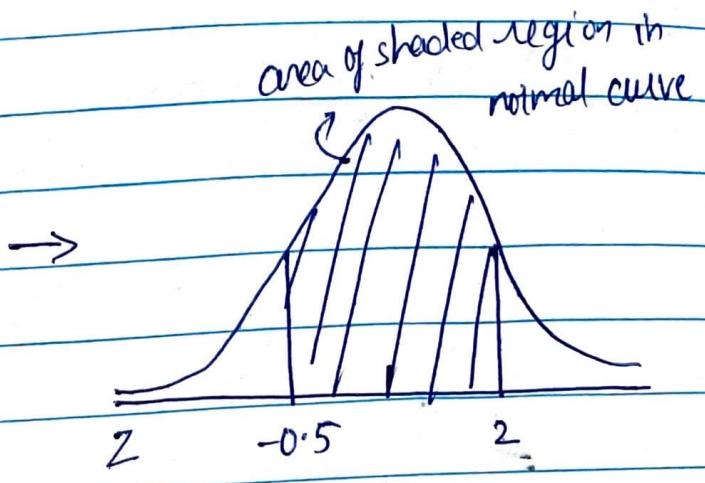
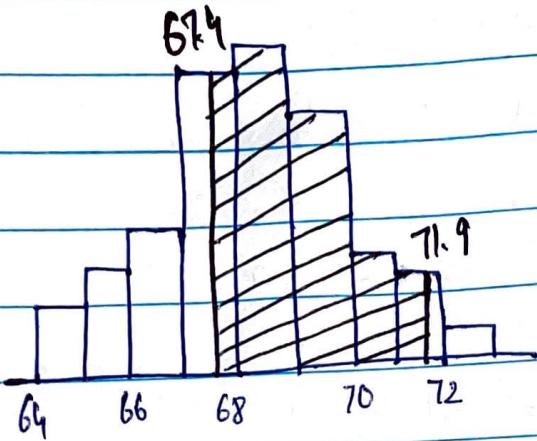
Q) What % data is b/w 67.4 & 71.9, when  $\bar{x}=68.3$  &  $s=1.8$

A) ① Standardize data

$$Z = \frac{67.4 - 68.3}{1.8} = -0.5$$

$$Z = \frac{71.9 - 68.3}{1.8} = 2$$

② Mark the area under normal curve



③ Rewrite this graph as = (all of area to left of 2) - (all of area left of -1.5)



④ Find values on table or software

$$\Rightarrow 97.7\% - 30.9\% = 66.8\%$$

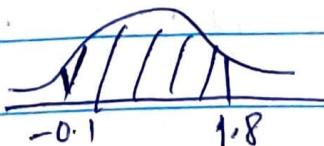
⑤ if  $\rightarrow 46\%$  &  $1.8 \rightarrow 96.4\%$

$$\bar{x} = 68.3$$

$$s = 1.8$$

$$z = \frac{68.4 - 68.3}{1.8} \approx -0.1$$

$$z = \frac{71.5 - 68.3}{1.8} \approx 1.8$$



$$\Rightarrow 96.4 - 46$$

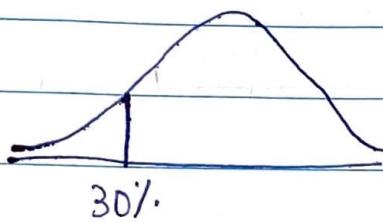
$$\Rightarrow 50.4$$

## Computing percentiles with normal approximation :-

Q) What is the 30<sup>th</sup> percentile of given data?

↓ where 30% of data falls below

usually,  $Z_1 = -0.52$



$$Z = \frac{\text{data} - \bar{x}}{s}$$

$$\Rightarrow Zs + \bar{x} = \text{data}$$

(or)

as z-score tells us how many standard deviations above or below average, so as  $Z = -0.52$ , it is 0.52 standard deviation below average of  $\bar{x}$ , so

$$\begin{aligned}\text{data} &= \bar{x} + Zs \\ &= 68.3 - (0.52 \times 1.8) \\ &= 67.4\end{aligned}$$

## Binomial setting and Co-efficient :-

Revisiting,  $P(\text{newborn is a girl}) = 49\%$ .

What are chances that if we look at 3 newborns, we have 2 girls?

$$\begin{aligned}P(2 \text{ of } 3 \text{ girls}) &= P(GGB \text{ or } GBG \text{ or } BGG) \\ &= P(GGB) + P(GBG) + P(BGG) \\ &= P(G)P(G)P(B) + \dots + P(B)P(G)P(G)\end{aligned}$$

ways we can arrange  $\cong (3) \times \underbrace{(0.49 \times 0.49 \times 0.51)}_{\text{same for all 3 terms}}$

2G and 1B

=

same for all 3 terms. This is binomial setting.

binomial setting:

$n$  = independent repetitions

each of these have 2 outcomes (success/failure)

$P(\text{success})$  = same in each experiment.

⇒ Now what is  $P(2 \text{ out of } 5 \text{ are girls})$ ?

If we try to do it as before, enumerating possibilities, the number increases as  $n$  does. In this case, it is 10 ways.

To ease if we use

$$\frac{n!}{k!(n-k)!}$$

binomial coefficient

ways to arrange "k" successes out of "n" ~~possible~~ independent repetitions.

∴ ways to arrange 2 girls out of 5 newborns,

$$\frac{5!}{2!3!} = \frac{5 \times 4 \times 3 \times 2 \times 1}{2 \times 1 \times 3 \times 2 \times 1} = \underline{\underline{10 \text{ possibilities}}}$$

Binomial formula:- (Binomial probability)

Applying this coefficient in binomial setting gives

$$P(k \text{ success in } n \text{ experiments}) = \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k}$$

probability of having a particular pattern of  $k$  success &  $(n-k)$  failures.

(Q) Play an online game 10 times. Each time there are 3 possible outcomes,  
 $P(\text{win big}) = 10\%$ ;  $P(\text{win small}) = 20\%$ ;  $P(\text{win nothing}) = 70\%$ .  
 $P(2 \text{ small prizes})?$

(A)  $P(\text{win big}) = 10\%$ . Success = win small prize  
 $P(\text{win small}) = 20\%$ . Failure = win big prize or win nothing  
 $P(\text{win nothing}) = 70\%$ .

$$\begin{aligned} P(2 \text{ success in } 10 \text{ exp}) &= \frac{10!}{2! 8!} (0.2)^2 (0.7)^8 \\ &= \frac{(10 \times 9)}{2 \times 1} \times 0.04 \times (0.7)^8 \\ &= 0.3019 \\ &= \underline{\underline{30.19\%}} \end{aligned}$$

## Random variables & Probability histograms :-

The outcomes of 'n' experiments is due to chance, so the num of success we see 'k' is random.

$X$  = 'num of successes' is a random variable.

$P(X=2) = 30.19\%$  in binomial distribution.

~~work~~

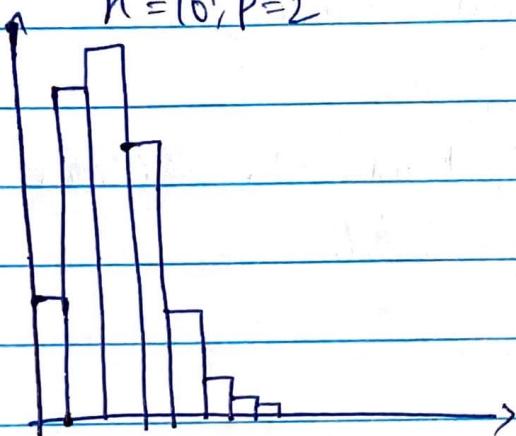
We can visualize the probabilities of various outcomes of  $X$  with probability histogram.

Height of bar =  $P(\text{Corresponding outcome})$

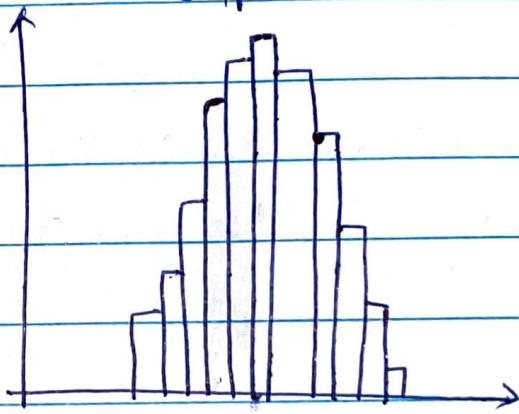
Normal approximation to the binomial :-

As number of experiment,  $n$ , gets larger, probability histogram of binomial distribution looks more similar to normal curve.

$$n=10, p=2$$



$$n=50, p=2$$



This suggest we can use normal approximation to compute binomial probability, subtract off  $\boxed{np}$

$$\boxed{\sqrt{np(1-p)}}$$

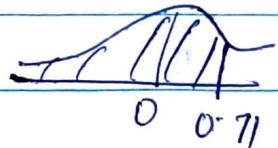
① Same question,  $n=50$ ,  $P(\text{at most } 2 \text{ small})$

$$p=0.2$$

Standardize, 
$$\frac{12 - np}{\sqrt{np(1-p)}} = \frac{12 - (10)}{\sqrt{50 \cdot 0.2 \cdot 0.8}} = \frac{2}{2.83} = 0.71$$

So, we have to find area under curve to left of 0.71

which is roughly 76%.



- ⇒ In the setting of survey, poll is a simple random sample, that is sampling without replacement. So this is not binomial setting, because 'p' changes after a subject has been removed. But if population is much larger then it has no effect.
- ⇒ A simple random sample selects subjects without replacement & each subject has equal chance of being selected.

① Coin is tossed 400 times. Chance of getting more than 210 tails?

(A)  $n = 400$

$k = 210$

$(n-k) = 190$

$p = 0.2$

$$\frac{400!}{210! 190!} (0.2)^{210} (0.8)^{190}$$

$$Z = \frac{\frac{210-np}{\sigma}}{\sqrt{np(1-p)}} = \frac{210 - (400 \times 0.2)}{\sqrt{(400 \times 0.2)(1-0.2)}} \\ = \frac{65}{4} = 16.25$$

⑥ Coin is tossed 6 times. Chances of getting 2 tails in each of first 3 and last 3?

Ⓐ n=6      TTH    TTH  
 $p=0.2$       THT    THT  
 $k=2$       HTT    HTT,  
 $\frac{6!}{2!4!} = \frac{15}{60}$  possibilities

$P(2 \text{ tails in } 1/3 \text{ and } 3/6)$  =  $P(2 \text{ tails in first } 3)$  and  $P(2 \text{ tails in last } 3)$

$$= \left( \frac{3!}{2!1!} (0.2)^2 (0.8)^1 \right) \left( \frac{3!}{2!1!} \right)$$

⑦ MCA with 5 ques, 4 possible answers, 1 correct, chance of getting 2 correct?

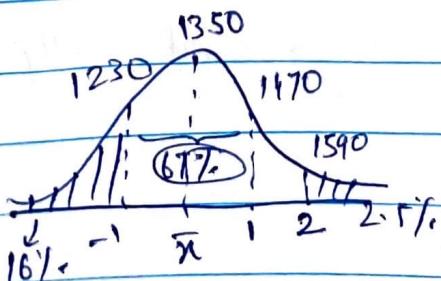
Ⓐ  $P(\text{success}) = 1/4 = 0.25$       ways =  $\frac{5!}{2!3!}$

$$\Rightarrow P(\text{wrong}) = 0.75$$

$$P(2 \text{ correct}) = \left(\frac{1}{4} \times \frac{1}{4}\right) = \left(\frac{1}{4}\right)^2$$

$$P(3 \text{ wrong}) = \left(\frac{3}{4}\right)^3$$

⑧  $\bar{x} = 1350$ ,  $s = 120$ , what % score below 1230 & what score is needed to be in top 2.5%.



$$z = \frac{1230 - 1350}{120} = -1$$

top 2.5%  $\Rightarrow$  2<sup>nd</sup> std deviation

$$\text{data} = z s + \bar{x} = 2(120) + 1350 \\ = \underline{\underline{1590}}$$

## Parameter & Statistic :-

- That in every statistic analysis, there are several histograms floating around. The probability histogram which generates the data. The histogram of observed data and the probability histogram of an estimate.
- We can estimate average height of population with a relatively small sample.

$\mu$  = average of population

$\sigma$  = standard deviation of population

Parameter → quantity of interest measured ~~in sample~~ about population

Statistic → quantity of interest measured in sample.

## Expected value & Standard error :-

If we draw a random sample from population, then we expect their size to be around the average population  $\mu$  give or take about  $\pm \sigma$ .

⇒ The expected value of one random draw is population average  $\mu$ .

⇒ Average of  $n$  draws,  $\bar{x}_n$  is  $E(\bar{x}_n) = \mu$

⇒ Since  $n$  is random,  $\bar{x}_n$  is also random, so not exactly is  $\bar{x}_n = \mu$ .

How far off from  $\mu$  is  $\bar{x}_n$ .

The standard error of a statistic tells roughly how far off the statistic will be from its expected value.

SE of a statistic  
SE plays same role as  $\sigma$  for random observation draws at random.

square root law,

$$SE(\bar{x}_n) = \frac{\sigma}{\sqrt{n}}$$

① It shows that SE becomes smaller if we use large samples  
we can use this formula for calculating SE.

② It does not depend on population size but on sample size.

$$\Rightarrow \text{Sum of } n \text{ draws} = S_n = n\bar{x}_n$$

$$\Rightarrow \text{Expected value of sum} = E(S_n) = n\mu$$

$$\text{Standard error of sum} = \sqrt{n}\sigma$$

Let's revisit what % of population approve way of P.M?

But the approval % is an average.

% of likely voters

→ Population falls into either of two categories. Put '1' on one who approves & '0' doesn't

Let there be 5 voters, 10010

sum of labels = 2 which is no. of people who approved.

∴ % of likely voters who approve is % of 1s among the labels.

So, in a sample of  $n$ ,

the number of voters in sample who approve is  $s_n$

$$\% \text{ approval} \Rightarrow \frac{s_n}{n} \times 100\% = \bar{x}_n \times 100\%$$

Expected value of  $s_n = N \times 100\%$ .

$$SE \text{ of } s_n = \frac{\sigma}{\sqrt{n}} \times 100\%$$

All these formulas are true even when data is simulated, generated according to histogram.

If a random variable  $X$  that is simulated has  $K$  possible outcomes  $x_1, x_2, \dots, x_K$ , then

$$\mu = \sum_{i=1}^K x_i P(X=x_i) \quad \sigma^2 = \sum_{i=1}^K (x_i - \mu)^2 P(X=x_i)$$

The square root law :-

⑥ Toss a coin 100 times. How many tails we can see? give or false.

(A) 1 = tails }  $P(0) = P(1) = \frac{1}{2}$   
 0 = heads }

no. of tails = sum of 100 draws

$$\Rightarrow E(S_n) = 100 \times \mu \quad \mu = \text{sum} = 0(P(0)) + 1(P(1)) = \frac{1}{2} \\ = 50$$

give or false,

$$SE(\text{sum}) = \sqrt{100} \sigma = 10 \times \frac{1}{2} = 5$$

$$\sigma^2 = (0 - \frac{1}{2})^2 \frac{1}{2} + (1 - \frac{1}{2})^2 \frac{1}{2} = \frac{1}{4}$$

so tails, give or false 5.

$$SE(\text{percentage}) = \frac{5}{\sqrt{100}} = 5\%$$

## The sampling distribution :-

Toss a coin 100 times. Possible tails range from 0 to 100  
How likely is each outcome?

This is a binomial distribution,  $n=100$  &  $p=1/2$  because if we call a coin landing success then no. of tails = no. of successes in 100 exp.

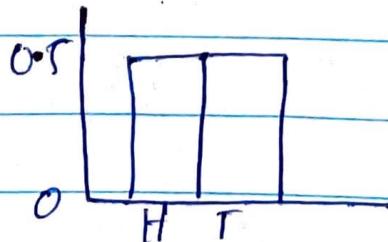
So, if statistic interest of  $S_n = \text{no. of tails}$ , then  $S_n$  is a random variable whose probability histogram is given by binomial distribution. This is called sampling distribution of statistic  $S_n$ .

Provides more details about chance properties

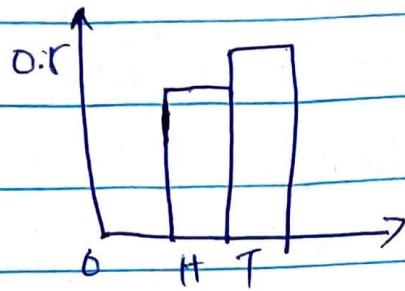
## 3 histograms :-

The chance process of tossing a coin 100 times with 3 different histograms

① Probability histogram  $\rightarrow$  theoretical



② Actual experimented histogram, let's say 53T & 47H



③ Sampling distribution : if we are interested in a statistic such as the no. of tails in 100 tosses, then we get this which gives Sampling distribution of that statistic

Law of large numbers :-

The square root law says that SE of sample mean goes to zero as sample size increases. So, the sample mean will likely to be close to its expected value  $\mu$  if the sample size is large. This is law of large numbers

→ Applies only for averages & %.

→ Sampling with replacement.

## Central limit theorem:-

Let's go back to the online game.  $n=10$ ,  $p=0.2$  (small win).

We gamble  $n$  times & found  $X$  has binomial dist.

When we sample with replacement, and  $n$  gets large, then the sampling distribution of sample average / sum % approximately follows normal curve.

To standardize, subtract expected value off statistic, then divide by SE.

$$\textcircled{Q} \quad \frac{\text{POP}}{\text{BS}} = 10,000 \quad ; \quad n = 100$$

$l = 6000 \quad \therefore l\% \text{ will be}$

$$0 = 4000$$

$$\textcircled{A} \quad n = 100 \quad \sigma_l = \sigma_0 = 0.5$$

$$P(l) = \frac{0.6}{100} = 0.6$$

$$P(0) = 0.4$$

$$E(S) = 100 \times 11 = 100 \times (0 \times P(0) + 1 \times P(l)) = 100 \times 0.6 = 60$$

$$SE \% = \frac{0.6}{\sqrt{100}} = \frac{6}{10 \times 10} = 60\%$$

$$SE = \sqrt{100} \sigma = 10 \times 0.5 = \underline{\underline{5}}$$

② 100 pledges, each pledge is equally likely to be \$10, \$50 or \$100.

$\sigma_{10} = \sigma_{50} = \sigma_{100} = \$37$ . Expected value of sum of 100 pledges?

Ⓐ expected value of sum =  $n\mu$

$$= 100 \times 4$$

$$\mu = \sum_{\substack{i=1 \\ i \neq 0}}^k x_i P(X=x_i) = \frac{\$10 + \$50 + \$100}{3} = \$53.33$$

$$\therefore \underline{E(S_n) = \$53.33}$$

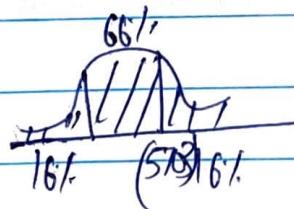
③ Same question. Chance that 100 pledges more than \$5700

Ⓐ  $SE = \sqrt{n}\sigma = \sqrt{100} \times 6 = \$370$

$$\frac{53.33}{370}$$

\$5703 - \$5700 is one standard deviation away,

so it's 16%.

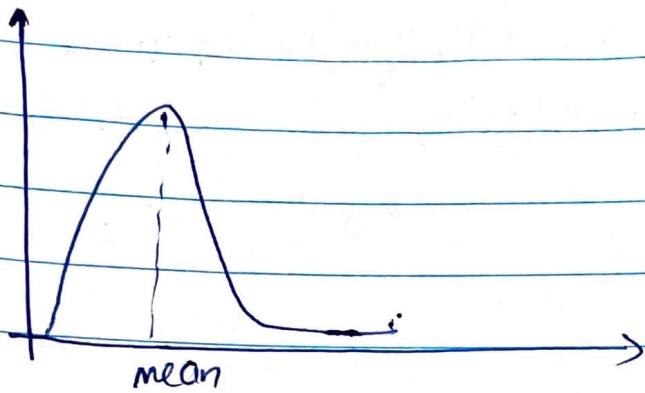


(Contd CLT):

It shows the statistic has normal distribution no matter what population histogram is.

Let  $\sigma = \$38,000$  in an income graph skewed right.

$$\mu = \$67,000$$



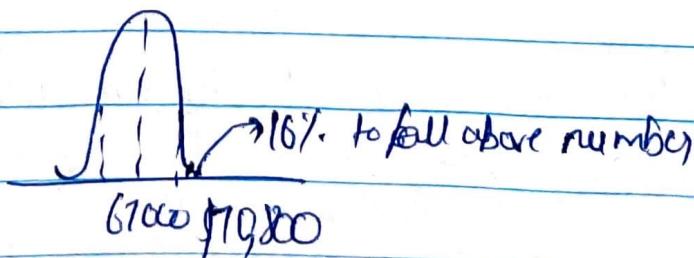
So, if we same "n" incomes are random, then

$\bar{X}_n$  = Sample average follows normal curve. To do normal approximation, we do subtract off the expected value of statistic which in this case is average of all incomes  $E(\bar{X}_n) = \mu = \$67000$  and then divide by  $SE(\bar{X}_n) = \frac{\sigma}{\sqrt{n}} = \frac{\$38000}{\sqrt{n}}$

① Let take 100 incomes, if  $n=100 \Rightarrow \sqrt{n}=10$

$$\therefore SE(\bar{X}_{10}) = \$3800$$

The CLT says, the sample mean follows normal curve centered at \$67000 with  $SE = \$3800$



① we did 'n' gamble and we had  $X = \text{no of small prizes}$ . Since we use labels,  $1 = \text{everything of small prize, everything else} = 0$

$\Rightarrow X = \text{no. of small prizes} = \text{sum of labels}$

Since we are looking at sum, we use CLT.

$\Rightarrow \mu = p \text{ & } \sigma = \sqrt{p(1-p)}$  for 1 gamble

$\Rightarrow \mu = np \text{ & } \sigma = \sqrt{np(1-p)}$