

Statistics

- Descriptive → you have All the Data & you want to Describe it
- Inferential → You have Portion of Data & you want to Draw conclusion.

Central tendency of the Dataset.

- Arithmetic Mean → $\text{sum}(\text{All the Numbers}) / \# \text{Numbers}$
- Median → middle element In the sorted data
- Mode → Most frequently occurring digit

Sample and population

$$\text{Population Mean} = \mu = \frac{\sum_{n=1}^N x_n}{N} \quad \vdots \quad \text{Sample Mean} = \bar{x} = \frac{\sum_{n=1}^n x_n}{n}$$

only the Notational Difference.

Measure of Dispersion

→ variance

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

For the population

For sample

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

But when we calculate the variance for sample we divide it by $(n-1)$ To get unbiased estimate

$$s^2_{n-1} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n-1)} \leftarrow \text{variance for sample.}$$

Standard deviation

$$\sigma = \sqrt{\sigma^2} = \sqrt{\text{variance.}}$$

Random Variable

- It's Function that maps Random processes To some value.

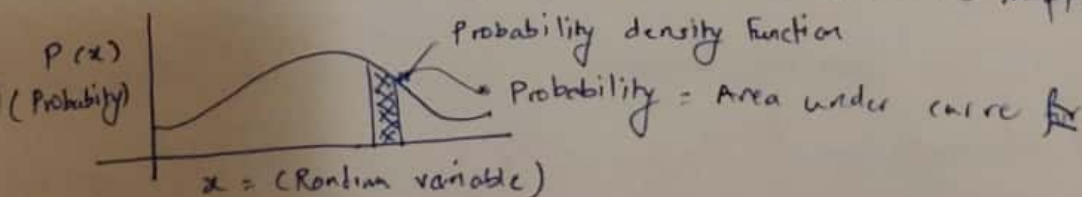
$$X = \begin{cases} 1, & \text{If it's Rain Tomorrow} \\ 0, & \text{If it's Not Rain} \end{cases}$$

mapping

$X \rightarrow$ can be discrete (Probability distribution)
can be continuous. (Probability density function)

Probability Density Function

- In case of continuous Random variable we have mapping function



Binomial Distribution

DATE -

Page - 2

→ Experiment consists of n Independent Trials. with Two mutually exclusive outcomes

→ so probability of success $= (P)$ then Prob. of failures $(1-P)$

Now you want to determine the probability of the success for n trials

X = Discrete random variable. = Number of success in n trials.

$$f(x) = {}^nC_x (P)^x (1-P)^{(n-x)}$$

The distribution you get is called as a binomial distribution

Expected Value $[E(X)]$

→ It is nothing but different representation of the mean.

→ Here you know the frequency of event's occurrence in the discrete case

→ and in continuous case you know the probability density function.

$$E(X) = \sum_{i=1}^n x_i p_i \rightarrow x_i = \text{set of outcomes of exp.}$$

p_i = Associated probabilities.

$$E(X) = \int x \cdot f(x) \cdot dx \quad \text{where } f(x) = \text{Probability density function}$$

x = Random variable at which you want to calculate the Prob.

Expected value for Binomial Distribution

$$E(X) = (\text{Probability of success}) \cdot (\# \text{ for which I want to calculate Exp value})$$

$$E(X) = n \cdot P$$

Poisson Process → (continuous Representation of Binomial Dist.)

X = # of car pass in an hour.

$E(X)$ = Let say you start to measure that λ no. of car passed in hour. We can model this experiment as a binomial distribution.

where n = number of minute

P = Prob. probability of car pass in a minute

$$E(X) = \lambda = n \cdot P$$

So now I have the Expected value so let calculate Prob. for 15 car per minute

$$P(X=k) = \left({}^{60}C_k \right) \cdot \left(\frac{\lambda}{60} \right)^k \left(1 - \frac{\lambda}{60} \right)^{60-k} \text{ — Binomial Dist}$$

As the $n \rightarrow \infty$ which is so above poisson distribution

and in Poisson Dist we measure Prob. for Interval and not for trial.

$$\lim_{x \rightarrow \infty} \left(1 + \frac{a}{x}\right)^x = e^a$$

We know from Binomial Distribution that

$$\lim_{n \rightarrow \infty} \frac{n!}{(n-k)! k!} P(X=k) = {}^n C_k \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} \quad \text{--- where } \lambda = \text{rate of success.}$$

As $n \rightarrow \infty$

$$P(X=k) \lim_{n \rightarrow \infty} = \frac{(n) \cdot (n-1) \cdot \dots \cdot (n-k+1)}{k!} \cdot \frac{\lambda^k}{k!} \cdot \left(1 - \frac{\lambda}{n}\right)^n \cdot \left(1 - \frac{\lambda}{n}\right)^{-k}$$

$$P(X=k) \lim_{n \rightarrow \infty} = \frac{n \cdot (n-1) \cdot \dots \cdot (n-k+1)}{n^k} \cdot \frac{\lambda^k}{k!} \cdot \left(1 - \frac{\lambda}{n}\right)^n \cdot \left(1 - \frac{\lambda}{n}\right)^{-k}$$

\downarrow $\frac{n \cdot (n-1) \cdot \dots \cdot (n-k+1)}{n^k} \rightarrow 1$
 \downarrow $\frac{\lambda^k}{k!}$ constant
 \downarrow $\left(1 - \frac{\lambda}{n}\right)^n \rightarrow e^{-\lambda}$
 \downarrow $\left(1 - \frac{\lambda}{n}\right)^{-k} \rightarrow 1$ As $n \rightarrow \infty$, $\lambda = \text{constant}$

$$P(X=k) \lim_{n \rightarrow \infty} = \frac{\lambda^k}{k!} e^{-\lambda}$$

The probability of car pass is now trial independent in a way its now continuous distribution.

So Poisson Distribution is special approximation of Binomial Dist. where No. of trials are infinite.

• Normal Distribution (Gaussian Distribution)

$$P(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

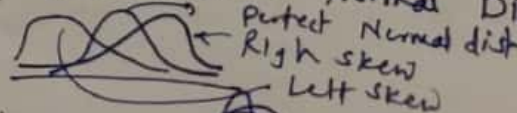
→ As opposed to the Binomial dist. ~ Poisson dist. The Normal dist. is continuous func.

→ Normal distribution Defined by The mean & standard deviation

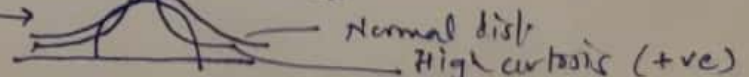
• Central limit Theorem

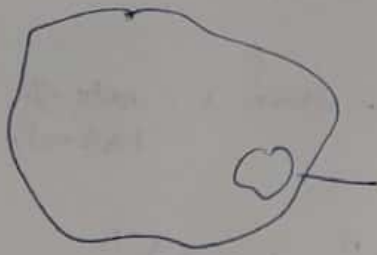
→ So Let's Assume you draw samples from identical same distribution with Re some sample size. The central limit theorem states that As your # of sample increases the (distribution) If you find the distribution of sample sum/mean It will be close to Normal Dist. & As $(n \rightarrow \infty)$ Dist. becomes pure Normal Dist.
 → As sample size & std. dev. of distribution dec. \rightarrow Dist. becomes pure Normal Dist.
 # → Samples → ∞ Sampling distribution → Normal Dist.

→ skew In Normal Distribution



→ Kurtosis In Normal Distribution





Population

Sample \rightarrow Sample mean
 \rightarrow sample variance.

\downarrow Now if you draw Infinite No. of Samples

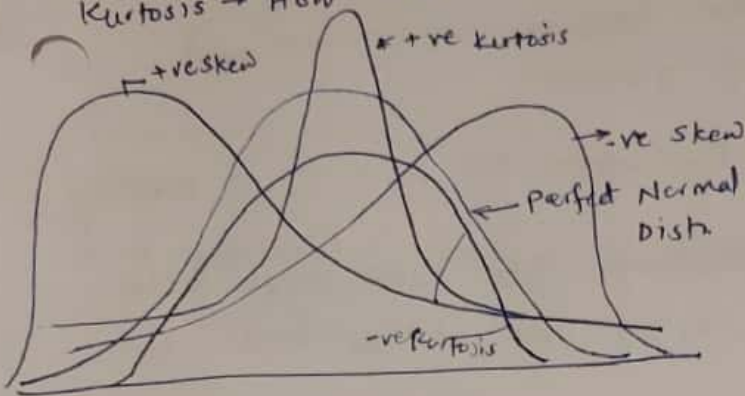
The mean will be close to population mean however standard deviation of Sample mean distribution

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

$\sigma_{\bar{x}}$ \leftarrow std for sample
 σ \leftarrow std for sample dist
 \sqrt{n} \leftarrow sample size.

Sampling Distribution of sample mean.

Skew \rightarrow How the distribution differ from Normal Dis horizontally
Kurtosis \rightarrow How the distribution differ from Normal Distribution vertically



\rightarrow CLT doesn't work if sample size is equal to 1

\rightarrow As the sample size $\rightarrow \infty$ the distribution looks like Normal Dist.

\rightarrow As the #sample $\rightarrow \infty$ the dist goes close to Normal Dist

\rightarrow As samples $\&$ #sample inc standard dev. decreases.

\rightarrow So Let's say you have a Non Normal some Random distribution with some mean & variance. the If you calculate Let's say you draw a sample of size n . and plot them on graph

\rightarrow The mean of the sample distribution \approx mean of original distribution

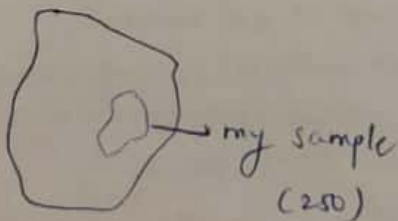
\rightarrow The standard deviation of $\sigma_{\bar{x}}^2$ for the sample distribution will be

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} \leftarrow \text{where } n \text{ is sample size}$$

So $\mu_{\bar{x}} = \mu$, $\mu_{\bar{x}} = \mu$ of original Dist $\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} \rightarrow$ where n sample size.

Margin of error

\rightarrow So Let's say for Pixel (100, 100) I'm getting 2000 values of Distances. How ever I can't use All 2000 values. So I selected ²⁵⁰ ~~100~~ Point values Randomly from 2000 values.



Values for
Pixel (100, 100)
(2000)

Conclusion:

95% chance That a Random \bar{x} is within $2\sigma_{\bar{x}}$ of Total population.

So we are 95% sure that 68% to 91% of points are Reliable.

to shrink the % gap we have to draw more samples.

- We have total 2000 callbacks of depth data
- Let's say ground Truth is 1000mm And we consider the point is Reliable If distance value is betn ± 5 mm.
- Now we draw 50 samples Randomly from above for fixed $(1,1)$ Randomly from above 2000 values
- We don't care about The Type of Distribution Those 2000 ^{Population} Sample follows. coz we know If we draw sample of size > 2 and plot the mean it will follow Normal Dist. ^{enough}
- So By central limit Theorem we are Rest Assure that The Distribution of Sample mean will be Normal Dist with some mean & standard dev
- Let say In our Sample of 50 we got 40 Reliable points & 10 unreliable points
- Now our task is to come up with the 95% Interval such that If we draw x points we are sure that some proportion of those x points are Reliable.
- As per central limit theorem we know

Population mean = mean of Sample distribution \approx sample mean.

Variance of Sample dist = $\frac{\text{variance of population}}{\text{sample size.}}$

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n}$$

For example

Sample Size = 50

Reliable point = 40

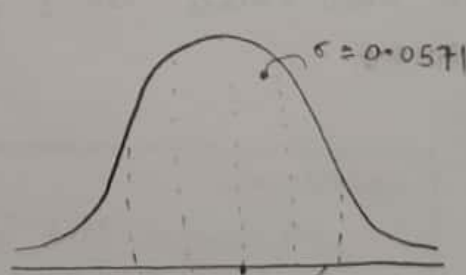
unreliable pts = 10

sample mean = 0.8

$$= \frac{40 \times 1 + 10 \times 0}{50}$$

$$\text{sample std dev} = \sqrt{\frac{(40 \times (1-0.8)^2 + 10 \times (0-0.8)^2)}{49}}$$

$$= 0.404 \rightarrow \text{std for Sample}$$



consist 95% Area.

$$\text{std for Sample Dist} = \frac{\text{std for sample}}{\sqrt{\text{sample size}}}$$

$$= \frac{0.404}{\sqrt{50}}$$

$$= 0.0571$$

Let's say we want 95% Confidence either side so it 2 σ Level
 \rightarrow confident 95% of chance that \bar{x} is within 0.8 Valid points is

\rightarrow Lower limit

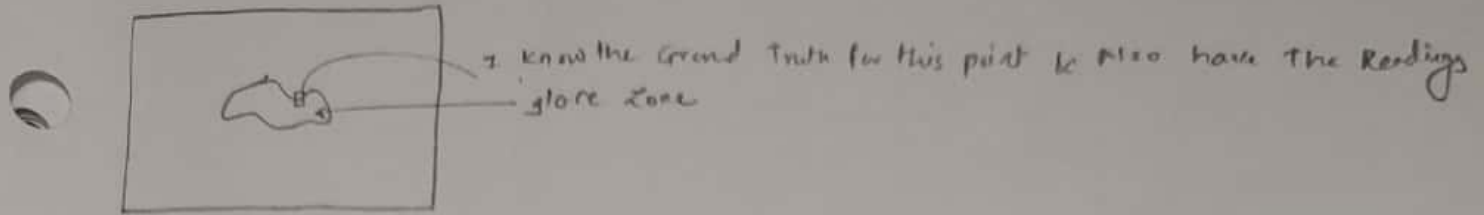
$$= 0.8 - 2 \cdot (0.0571)$$

$$= 0.6857$$

upper limit

$$= 0.8 + 2 \cdot (0.0571)$$

$$= 0.9142$$



Image

$H_0 \rightarrow$ Glare has no effect on the depth value for that point

$H_1 \rightarrow$ Glare has an effect

Let's Assume

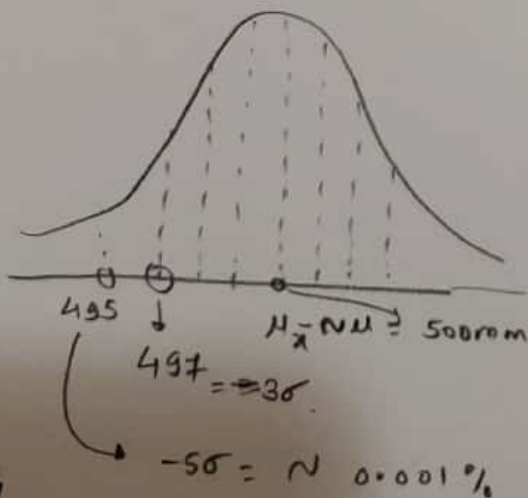
- \rightarrow We have 100 reading for the point under consideration
- \rightarrow We know the Ground Truth depth is 500mm
- \rightarrow For the 100 reading of point on glare boundary the mean is 495mm & $\sigma = 10$ mm

Now the process

- Calculate Standard deviation for Sampling distribution

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{\text{sample size}}} = \frac{10}{\sqrt{100}} = 1 \text{ mm}$$

- How many standard deviation away from Ground truth i.e. 500mm the sample mean 495mm is. so it's just like calculating the Z score.

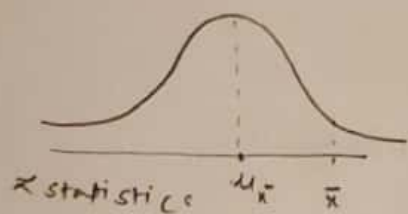


So the hypothesis that Glare has No effect is Rejected
 Because it has an chance of less than 0.01%. So

Z-Statistics vs T-Statistics

PAGE-5-B

We use The T-statistics to calculate The calculation steps are same as Z statistics. But even in case of T-statistics the distribution will look like T distribution.



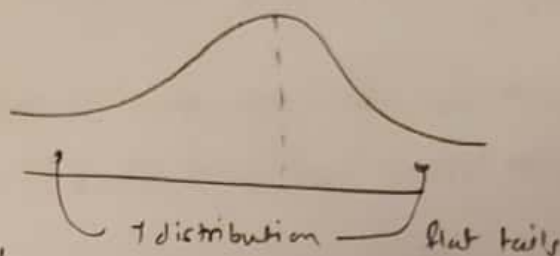
$$Z\text{-score} = \frac{\bar{x} - \mu}{\sigma_{\bar{x}}}$$

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

In case

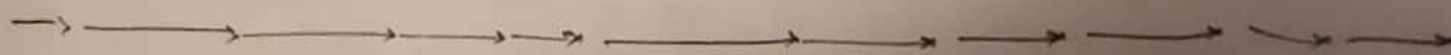
n is very small

we use T-score instead of Z-score



⇒ Type-I error

Rejecting the Null hypothesis even if it's true.



so we have

Sample size = 10

mean of sample = 17.17

Sample - std = 2.98

Now we want to find the

Zone with confidence level 95%

i.e. Z score = 2

$$-2.262 < Z < 2.262$$

We know $Z = \frac{\bar{x} - \mu}{\sigma_{\bar{x}}}$

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

σ ← sample std
√n ← sample size

$$\bar{x} = \bar{x} \rightarrow \text{sample mean}$$

$$Z = \frac{17.17 - \mu}{\frac{2.98}{\sqrt{10}}}$$

$$-2.262 < Z < +2.262$$

$$-2.262 (0.942) < 17.17 - \mu < 2.262 (0.942)$$

$$-2.13 < 17.17 - \mu < 2.13$$

$$2.13 > \mu - 17.17 > -2.13$$

$$19.3 > \mu > 15.04$$

95% chance that sample mean is betⁿ this Range.

Rule of thumb

$$n > 30$$

$$np > 5$$

$$n(1-p) > 5$$

Question 1 → Given a Distance value for a pixel with and without glare
Prove whether the glare have statistically significant effect on point's depth value

What we know →

① we have Sample of 2000 points for the Depth data of a pixel with & without glare value.

Our hypothesis

- H_0 → Glare has no effect on the depth value for that point
- H_1 → Glare has an effect on the depth value for that point

Our methodology

→ First calculate the ground Truth value

Ground Truth value = average (All point's without introduction of Glare)
for that particular Pixel

→ calculate the mean & standard deviation for the Glare sample

- collect the Pixel in the Glare boundary Region.
- get the value of distance for that Pixel. (sample > 200 points)
- Calculate the mean & standard deviation for the Glare point sample

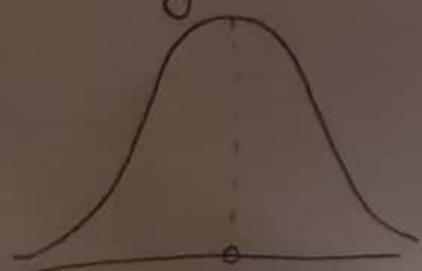
→ H_0 ⇒ Glare has no effect

→ H_1 ⇒ Glare has effect

Glare { $\mu_{\text{distance}} = \mu_{\text{true value}}$ (even with glare)
Points $\mu_{\text{distance}} \neq \mu_{\text{true value}}$

→ Assume H_0 is True

Sampling distribution



$\mu_{\bar{x}} = \mu = \text{True depth value}$
↓
mean for points w/o Glare

Sampling dist

Standard dev for sampling dist $\Rightarrow \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$

σ → std dev for population
 \sqrt{n} → sample size
 $\sigma_{\bar{x}} \approx \frac{s}{\sqrt{n}}$ → std of sample

So now we have

- ① Sample distribution mean $\mu_{\bar{x}} = \text{true depth value}$
- ② Sample → std. dev $\sigma_{\bar{x}} = \frac{s}{\sqrt{n}}$

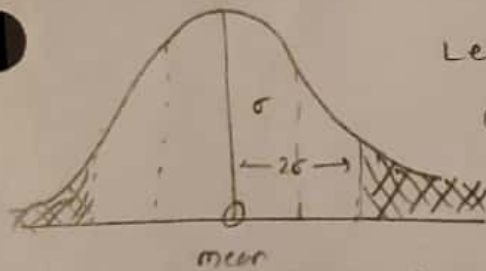
→ How many std. dev away from the sampling distribution mean is the mean of the glare sample & the what is the prob. of getting that Results away from sampling distribution mean.

→ Need to calculate the Z-score

$$Z = \frac{\text{(True depth value)} \text{ or } \text{(Sampling dist mean for the Pixl where there is No GLARE)}}{\text{(Glare depth value)} - \text{mean we are getting for the same Pixl where there is a (GLARE)}}$$

$$\text{Standard dev. of sampling distribution} = \sigma_{\bar{x}} = \frac{S}{\sqrt{n}} = \frac{\text{Standard dev. of Glare sample}}{\sqrt{\text{sample size}}}$$

→ once we get the Z value we are sure that the Result is Z standard deviation away from the mean. (True depth value)



Lets say our $Z=2$ then we are focusing on shaded area.

(sampling distribution for the Dataset w/o glare)

→ We know for $Z: 2\sigma$ The area under Normal curve where $K \text{ std } 2\sigma = 5\%$

→ so Based on the Z score we can say that

→ If the Null Hypothesis (H_0) is true then there is 5% chance of getting this value

→ so we can say that if we are confident 95% time for

Question No- 02

PAGE 7

→ Let's say you have set of 100 points distributed irregularly for the Particular Pixel

→ you can also know the True depth value for that pixel.

→ you consider the point as Reliable if its value is in some Threshold α

→ So Reliable point's Range is (true value $\pm \alpha$)

→ Based on this you divided your sample in Reliable and unreliable points.

→ Let's say P is the probability that point is Reliable
 $(1-P)$ → unreliable

→ Question → If we draw a sample of 100 point for particular pixel And then Base on some threshold 'T' from the mean of sample we decide whether the current point is reliable or not. Let say in sample we got $x\%$ of Reliable points. What is the 95% confidence Region that the Actual percentage of Reliable points in the population is close to sample mean.

Methodology: → [For $n = 100$, $(1-P) = 0.57$, $P = 0.43$, $n = 100$]

① First you have to calculate the mean & std. dev for a sample

$$\text{mean} \rightarrow \bar{x} = \frac{57 \cdot 0 + 43 \cdot 1}{100} = 0.43$$

std dev =

$$S^2 = (57 \cdot (0 - 0.43)^2 + 43 \cdot (1 - 0.43)^2) / 99$$

② We know mean of sampling dist of the sample mean $\mu_{\bar{x}} = \mu =$ of population

③ We know the std. dev of sample dist of sample mean = $\sigma_{\bar{x}}$
 $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$ $\sigma \approx \sqrt{S^2}$ → coz sample is from some dist

④ In statistic we know (current scenario)

Sample

$$\text{mean} = \bar{x} = \frac{0 \cdot q + 1 \cdot P}{n}$$

$$\text{std dev } S = \sqrt{\frac{P(1-\bar{x}) + \bar{x}(0-\bar{x})}{n-1}}$$

size → n

Sample dist sample mean

$$\mu_{\bar{x}}$$

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

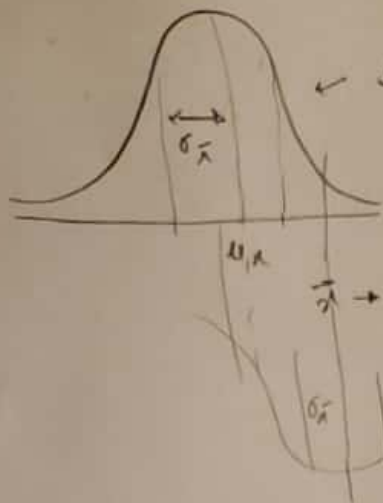
equal to

$$\sigma = S$$

Relationship betⁿ
 $\bar{x} \approx \mu$

Population
 μ

So we know value for \bar{x} , S , σ , $\sigma_{\bar{x}}$ and relation that $\mu_{\bar{x}} = \mu$ now Base on this we know to find



← So this is the distribution of our sampling means

$\bar{x} \rightarrow$ Lets assume \bar{x} for μ have So Based on \bar{x} we need To predict possible values for the μ

(5) So Based of It we can calculate the Z score.

$$Z = \frac{\bar{x} - \mu}{\sigma_{\bar{x}}}$$

We know α & σ & Range for Z value Based of normal dist table & we need to cal μ for 95% confidence

$$Z < \frac{\bar{x} - \mu}{\sigma_{\bar{x}}} < Z$$

$$Z = 2.58 \cdot \sigma_{\bar{x}}$$

$$-2.58 \sigma_{\bar{x}} < \bar{x} - \mu < 2.58 \sigma_{\bar{x}}$$

$$+2.58 \sigma_{\bar{x}} > \mu - \bar{x} > -2.58 \sigma_{\bar{x}}$$

$$+2.58 \sigma_{\bar{x}} + \bar{x} > \mu > -2.58 \sigma_{\bar{x}} + \bar{x}$$

Results: \rightarrow

Based of this results we can find the μ of an reliable point
Range of

Based of a sample.

Let's say we draw 100 points for some pixel & 100 points for some other pixel. Then we took a difference of mean for both samples. Now we want to say whether the points follow belong to different surface under consideration.

→ Let's say we have exact same position for the camera & we just change the surface under consideration does the sample data we get is enough to tell us the diff bet surface

→ our methodology →

For sample one we know
 mean → \bar{x}_1
 std dev = s_1 → $\sigma_{\bar{x}_1} = \frac{s_1}{\sqrt{n}}$
 for sample two we know
 mean → \bar{x}_2
 std dev → s_2 $\sigma_{\bar{x}_2} = \frac{s_2}{\sqrt{n}}$

Our sampling strategy is diff. Now we sample two diff values from two samples & subtract them so we get



② now we know

$$\mu_{\bar{x} - \bar{y}} = \mu_{\bar{x}} - \mu_{\bar{y}}$$

$$\sigma_{\bar{x} - \bar{y}}^2 = \sigma_{\bar{x}}^2 + \sigma_{\bar{y}}^2$$

$$\sigma_{\bar{x} - \bar{y}} = \sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}}$$

+value. →

for this case we know

x_1	x_2	$\mu_{\bar{x}_1 - \bar{x}_2}$	$\mu_{x_1 - x_2}$
s_1	s_2	$\sigma_{\bar{x}_1 - \bar{x}_2}$	σ

→ so problem here is given $\sigma_{\bar{x}_1 - \bar{x}_2}$ & $\bar{x}_1 - \bar{x}_2$ → calculate the 95% confidence interval for $\mu_{\bar{x}_1 - \bar{x}_2}$

③ Now we have to get the 95% confidence region $\sim 1.96 \sigma_{\bar{x}_1 - \bar{x}_2}$

confidence region = mean $\pm 1.96 \sigma_{\bar{x}_1 - \bar{x}_2}$
 of BP

So by other way = $(\bar{x}_1 - \bar{x}_2) \pm 1.96 (\sigma_{\bar{x}_1 - \bar{x}_2})$
 In terms of sample mean

→ The confidence interval will give the range in terms of the mean value
Now Based on the confidence interval we can write hypothesis

(3) $H_0: \mu_1 - \mu_2 = 0 \Rightarrow \mu_{\bar{x}_1} - \mu_{\bar{x}_2} \Rightarrow \mu_{\bar{x}_1 - \bar{x}_2} = 0$

Question No-4

given the data of a point determine whether or not point is in
Gore area

Breeman Science (Student's T-test)

once we get T-value $\rightarrow \sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}}$

→ we create a null Hypothesis with some confidence interval

Let's say $p=0.05$ i.e. 95% confidence

→ Now we need to find the degrees of freedom
 $= n+m-2$

→ we get the σ value for normal distribution

→ If the σ value is greater than our T value

we reject our Null Hypothesis (H_0)

→ we can do T test in spread sheet quite easily.

→ Population dist should be normal

→ similar variance

→ If the point is 20-30 data point we go for z test