

Unit 5 Stats

A/B testing, hypothesis tests, null hypothesis & alternative hypothesis, 1-way vs 2-way hypothesis testing, resampling, permutation test, p-values, t-tests.

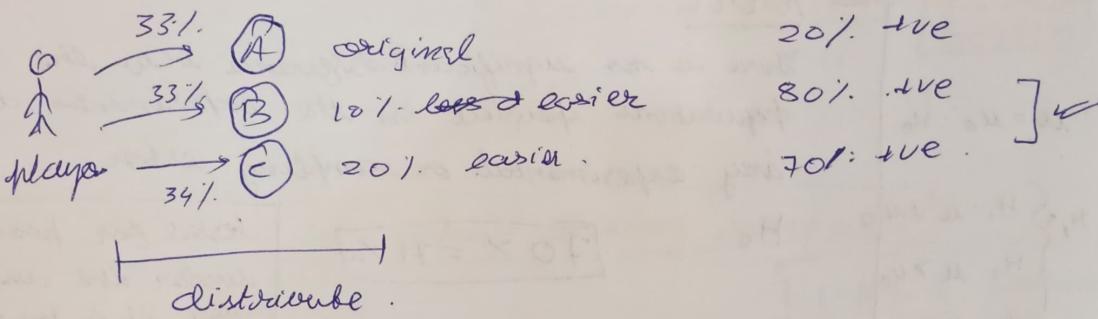
→ Question // Assumption.

A/B Testing:

is testing different versions at the same time & seeing which version gets the best response from players.

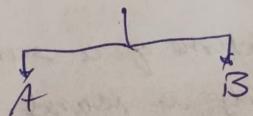
A version is compared with B version.

tests 1 variable
monitors one/two result.



① SUTVA (stable unit treatment value assumption)

→ the potential outcome on 1 unit should be unaffected by particular assignment of treatments to the other units.



A purchase independent of B is purchase

② Sampling:

→ represents a population.

distribution

→ sample distribution is representative of population.

conversion rate

Hypothesis Testing

the claims.

- It is a test to test some hypotheses about parent population from which sample is drawn.

class \rightarrow 100 no. % age 70%.

↳ 20 no. % age 71%.

→ We use hypothesis testing to check the difference b/w this 70% & 71%.

→ Is this 1% significant?

$\mu = \mu_0$: two tailed
$\mu > \mu_0$: right tailed
$\mu < \mu_0$: left tailed

population has (μ_0)

NULL

There is no significant difference b/w the populations specified in the experiments due to any experimental or sampling error.

$$\mu = \mu_0 : H_0$$

$$H_1 = \begin{cases} H_1: \mu \neq \mu_0 \\ H_1: \mu > \mu_0 \\ H_2: \mu < \mu_0 \end{cases}$$

$$H_0$$

$$70\% = 71\%$$

Tested for possible rejection under the assumption that it is true.

ALTERNATIVE

Observations are easily influenced by some random cause.

$$H_a \text{ or } H_1$$

$$70\% \neq 71\%$$

Eg: The mean lifetime of a sample of 31 glucometer produced by a company is found to be 58 months with SD of 10 months. Company claims that the mean lifetime of the glucometer is 60 months. Test the hypothesis @ 5% level of significance, whether the 58 months lifetime is acceptable or not?

dimmed by me mistake

70%
90%

This is a
Rasha.

$$n = 81$$

$$\lambda_{mean} = 58$$

$$SD = 10 (\sigma)$$

H_0

company. $\lambda = 60$

True

acceptance 58 or not?

Is the 2 months diff tested or not?

NULL

$$58 = 60.$$

ALTERNATIVE

$$58 \neq 60.$$

① Standard

Error Mean. (SEM)

$$\boxed{\frac{10}{\sqrt{81}}} = \frac{10}{9} = \boxed{1.11}$$

②

$$Z\text{-score} = \frac{\text{Difference}}{\text{SEM}} = \frac{|58 - 60|}{1.11} = \boxed{1.80}$$

Null & alternative contain differing endpoints

$1 - \beta \approx 1$ success

Null $\rightarrow = \text{always} \dots \geq \leq$

alternative $\rightarrow < \text{ or } >$

Errors in Sampling

$1 - \beta \Rightarrow$ reject null hypothesis when it is (H_0) is false.

Type 1 (

reject Null (H_0) when it is H_0 indeed true.

Type 2.

D.N.R
~~accept~~ Null (H_0) when H_1 is true
~~strong~~. (H_0 is false).

$$\alpha = P(\text{reject } H_0 | H_0 \text{ is true})$$

$$\beta = P(\text{accept } H_0 | H_1 \text{ is true})$$

Level of Significance: probability of taking risk of accepting rejecting a null hypothesis.

maximum α (ptable) risk.

Procedure for Hypotheses Testing

- ① set up hypotheses

↳ H_0 (Null)

↳ H_1 (Alternatives)

- ② set up significance level. 0.1 1 5 10

- ③ Test criteria

z test, < 30 students t test,
X test, F test

- ④ computation

- ⑤ decisions : TRUE / FALSE

Accept H_0

Reject H_0

H_0 is true.	correct	Type I error (α)
H_0 is false.	Type 2 error (β)	correct

- ② Set up significance level.

$$\alpha + \beta = 1$$

Significance level \rightarrow α = probability of occurrence (wrong decision to reject)

$$\boxed{\text{Significance level} + \text{confidence level} = 1}$$

Procedure for Hypotheses Testing.

- ① set up hypotheses
 - ↳ H_0 (NULL)
 - ↳ H_1 (ALTERNATIVES)
- ② set up significance level. 0.1 , 5 , 10
- ③ Test criteria :
 Z test, < 30 students t-test,
 X test, F-test
- ④ computation
- ⑤ decisions : TRUE / FALSE

	Accept H_0	Reject H_0
H_0 is true	correct	Type I error (α)
H_0 is false	Type II error (β)	correct

- ② Set up significance level.

$$\alpha + \beta = 1$$

Significance level \rightarrow α : probability of occurrence (wrong decision type I error)

$$\boxed{\begin{matrix} \text{Significance level} & + \text{confidence level} & = 1 \end{matrix}}$$

One-tailed & two-tailed:

$$Z = \frac{X - \mu_0}{\sigma / \sqrt{n}}$$

Eg: $H_0: \mu = 500 \text{ ml } (\mu)$

$H_1: \mu \neq 500 \text{ ml } [2 \text{ tailed test}]$

or $H_1: \mu < 500 \text{ ml } [\text{left test}] \quad \{ \text{One tailed.}$

or $H_1: \mu > 500 \text{ ml } [\text{right test}]$

When to use which H_1 ??

H_1

\neq

Eg: company claims that on average contains 100 mg of flour per kg. Test whether there is a strong evidence that it is greater than the claim.

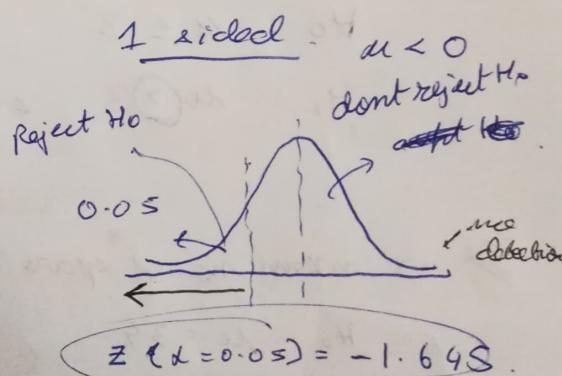
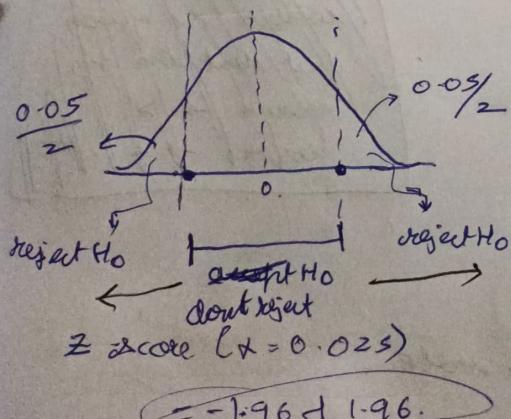
$H_0: \mu = 2$

$H_1: \mu > 2$

H_1
greater: $\mu > 2$
less: $\mu < 2$
differs: $\mu \neq 2$

Consider $\alpha = 0.05$.

2 tailed test: ($\mu \neq 0$)



(H_1)
based on alternative testing H_0 is either accepted or rejected
F.T.R

1-tailed:

- ① no detection of other side.
- ② directional power is gained.

Eg: Went to test if col students take less time than 5 years to graduate from col on average than.

$$H_0: \mu \leq 5 \quad (\text{position starting from})$$

$$H_1: \mu > 5 \quad (\text{where aiming?})$$

Eg: con test, 40% pass. we want to test if more than 40% pass that day.

$$H_0: \mu = 40 \quad (0.4)$$

$$H_1: \mu > 40 \quad (0.4) \alpha$$

If the H_0 is not given, I suppose you are testing if next connection more than 3 miles, then obviously μ or testing if you think it is less than or equal to 3.

This becomes Null hypothesis then.

Eg: testing mean speed is more than 3 miles.

$$H_0: \mu \leq 3$$

$$H_1: \mu > 3$$

Eg: 2) mean no. of years is 34

$$H_0: \mu = 33$$

$$H_1: \mu \neq 34$$

2) at most 60% of americans vote.

$$H_0:$$

$$H_1: \mu$$

Practically:
whatever is given
just that is
more than $\rightarrow >$
less than $\rightarrow <$, etc.

Null:

Assume H_0 is true. (claim)

↓
reject H_0

↓
Fail to reject H_0

If you want to support a claim YOU MUST state it as H_1 (not H_0)

want to prove it? H_1

NOTE:

- ① The original claim could either be H_0 or H_1 ,
 ↳ depends where EQUALITY IS

Eg: mean of fluid is AT LEAST 12 oz in a can.

CLAIM:

$$\mu \geq 12$$

H_0

OPP:

$$\mu < 12$$

H_1

which is H_1 or H_0 ?

whichever there is equals ... put it clear.

Eg: 1 mean no of years americans work is 34

$$\mu = 34 \quad H_0$$

$$\mu \neq 34 \quad H_1$$

Eg: at most 60% vote in elections.

$$p \leq 0.6 \quad H_0$$

$$p > 0.6 \quad H_1$$

Eg: mean starting salary is at least \$10000 per year

$$\mu \geq 10000 \quad H_0$$

$$\mu < 10000 \quad H_1$$

Eg 5 21. get drunk

$$\mu = 0.29 \quad H_0$$

$$\mu \neq 0.29 \quad H_1$$

$$n = 0.89$$

$$n > 0.44$$

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

Eg 6 fewer than 5% of adults ride the bus.

$$\mu < .05 \quad H_1$$

$$\mu \geq 0.05 \quad H_0$$

Eg 7 mean no. is not more than 10

$$\text{mean } \mu \leq 10 \quad H_0$$

$$\mu > 10 \quad H_1$$

Eg 8 about half prefer to stay outside.

$$\mu = 0.5 \quad H_0$$

$$\mu \neq 0.5 \quad H_1$$

Eg 8 chance of developing breast cancer is under 11% of women

$$\mu < 0.11 \quad H_1$$

$$\mu \geq 0.11 \quad H_0$$

$$PC \times F(10) \approx 81\%$$

Eg 9. mean cost is more than 2k

$$\mu > 2k \quad H_1$$

$$\mu \leq 2k \quad H_0$$

Eg 10 $m = 273 \quad \gamma = 0$

$m = 63$ thin. $\gamma = 4$

Is there a good evidence to suggest more than 30% of
teen girls stay thin by smoke? $H_1 = ?$

$$H_{01} =$$

$$\mu > 0.3 \quad H_1$$

$$\mu \leq 0.3 \quad H_0$$

$$n = \frac{63}{273} \times 23 =$$

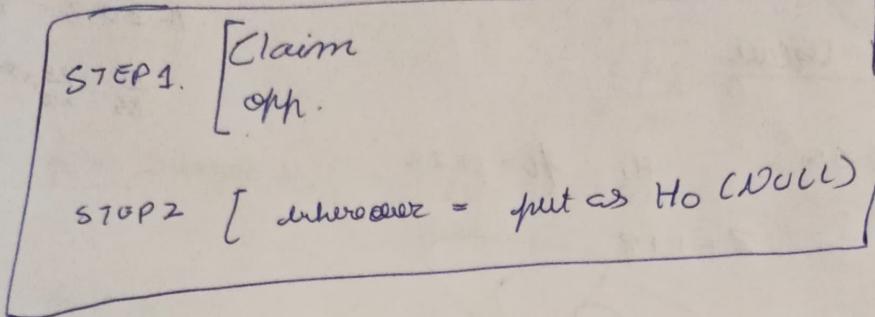
23

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

$$T = \frac{\bar{x} - \mu}{S / \sqrt{n}}$$

s = population S.D.

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$



Eg: sample of 706 found (61%) were A.C. claim:
most cases are male

$$\text{CLAIM: } \mu > 0.5 \Rightarrow H_1: \mu > 0.5$$

$$\text{opp: } \mu \leq 0.5 \Rightarrow H_0: \mu = 0.5$$

$$\hat{\mu} = 0.61 \quad \mu = 0.5 \quad z = 1 - \mu$$

$$Z = 5.84$$

\leftarrow UNUSUAL.

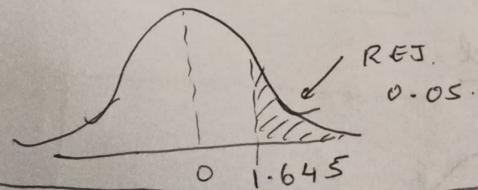
H_0 is rejected.

so H_1 is accepted.

HOW TO MAKE DECISION?

significance level:

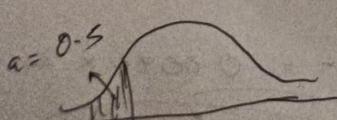
$$\alpha = 0.05$$



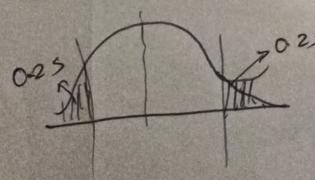
$$(H_1: \mu < 0.5)$$

$$(H_1: \mu = 0.5)$$

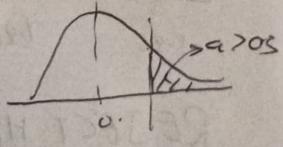
$$H_1: \mu > 0.5$$



left tail.



Two tail



right tail.

P-Value



REJECT H_0 IF P-value $\leq \alpha$

F.T.R.

H_0 IF P-value $> \alpha$

Eg:

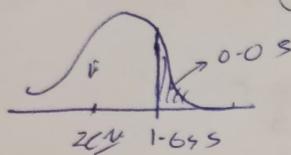
Find P-values

$$\textcircled{1} \quad \alpha = 0.05$$

$$H_1: \mu > 0.25$$

$$Z = 1.18$$

TRAP



F.T.R
 H_0 ~~Rejected~~

not enough evidence
to support the claim.

F.T.R

P-value

$$0.8810$$

$$Z = 1.18$$

test stat.

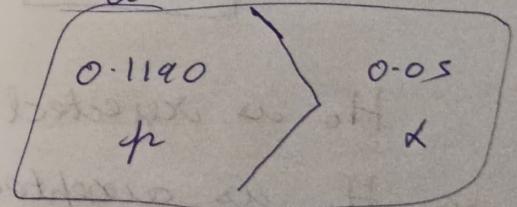
$$\alpha = ?$$

P-value

OOPS!!

table gives left area to Z.

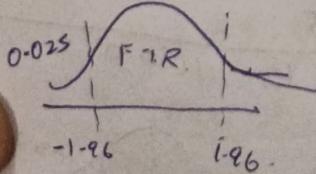
$$\alpha = 1 - 0.8810 = 0.1190$$



$$\textcircled{2} \quad \alpha = 0.05 \quad H_1: \mu \neq 25$$

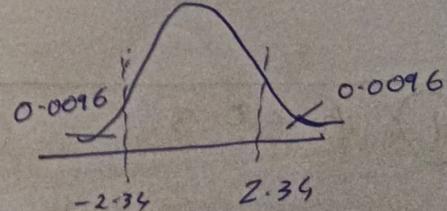
$$Z = 2.39$$

trap:



REJECT H_0 .

P-value-



$$\text{P-value} = 0.0096 \times 2$$

$$0.0192 < 0.05$$

reject H_0 .

P-Value



REJECT H_0 IF P-value $\leq \alpha$

F.T.R

H_0 IF P-value $> \alpha$

In a

that

class

OP

Eg:

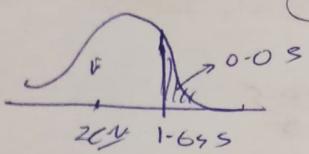
Find P-values

$$\textcircled{1} \quad \alpha = 0.05$$

$$H_1: \mu > 0.25$$

$$Z = 1.18$$

TRAD



H_0 F.T.R
Rejected

not enough evidence
to support the claim.

P-value

P-value

$$\begin{array}{ll} \mu < 0.2 & H_1 \\ \mu \geq 0.2 & H_0 \end{array}$$

$$\frac{11}{84} \quad \frac{20}{48} \times \frac{84}{120} = \frac{1}{6}$$

$$0.8810$$

$$Z = 1.18$$

tail int.

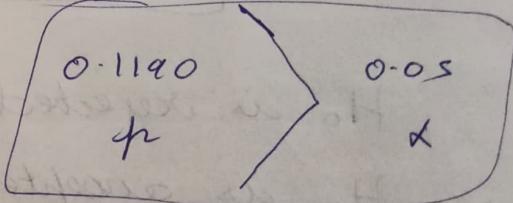
$$\begin{array}{l} \alpha = ? \\ \text{P-value} \end{array}$$

$$\alpha = 0.8810 \quad (\text{far})$$

OOPS!!

tools give left area to Z.

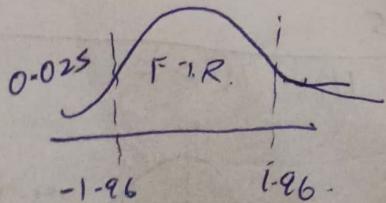
$$\alpha = 1 - 0.8810 = 0.1190$$



$$\textcircled{2} \quad \alpha = 0.05 \quad H_1: \mu \neq 0.25$$

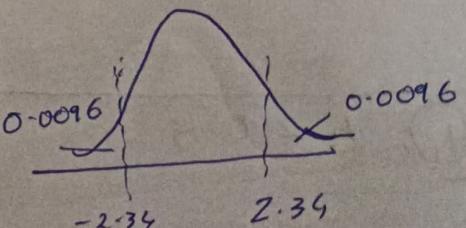
$$Z = 2.34$$

Trad:



REJECT H_0 .

P-value



$$\text{P-value} = 0.0096 \times 2$$

$$0.0192 < 0.05$$

reject H_0 .

In a sample of 300, it was found that most CEO were male. use significance level as 0.05

claim $\mu > 0.05 \Rightarrow H_1$ $\alpha = 0.05$.
opp $\mu \leq 0.05 \Rightarrow H_0$

$$Z = \frac{\hat{p} - p}{\sqrt{\frac{pq}{n}}} = 3.81$$

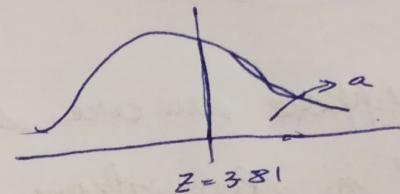
$$\hat{p} = \text{sample } p = \frac{183}{300} = 0.61$$

$n \geq$

$$p = 0.05$$

$$\alpha = 0.05$$

$$n = 300 \quad (\text{sample size})$$



accept H_1

support claim.

$\hat{p} \xrightarrow{\text{to}} \cancel{p} \times \text{reject } H_0$

reject H_0

Eg: previously organization has on average 4.5 random.

currently thinks the mean is higher is chosen

$$\text{mean} = 4.75, \sigma_{\text{sample}} = 2 \quad \alpha = 0.05$$

$$\mu > 4.5 \quad H_1$$

$$\mu \leq 4.5 \quad H_0$$

$$Z = \frac{4.75 - 4.5}{\frac{2}{\sqrt{15}}}$$

F.T.R

value.

About Pop Mean (σ known)

① random sample.

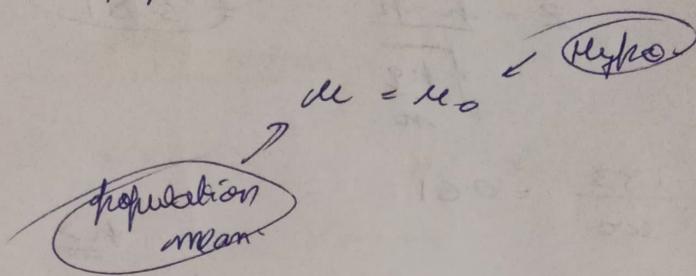
② σ is given.

③ $n > 30$ OR pop is normally distributed. \rightarrow ch 6

T-test

for population mean μ .

→ when sampling from normally distributed population.



Suppose we are drawing a simple random sample of n observations from normally dist. population.

To test $H_0 : \mu = \mu_0$

either Z-test or t-test

σ_{pop} is known

σ_{pop} is NOT known

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma_{\text{pop}}}{\sqrt{n}}}$$

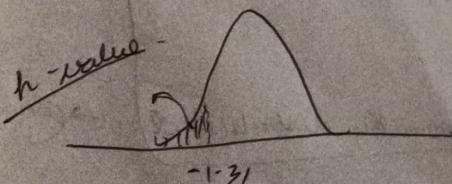
$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \quad S: \text{sample}$$

Standard error

estimated standard dev. of sampling distribution of sample mean.

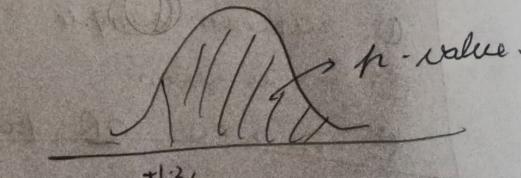
$$H_0 : \mu = \mu_0$$

$$\mu < \mu_0$$



$$t = -1.31$$

$$\mu > \mu_0$$



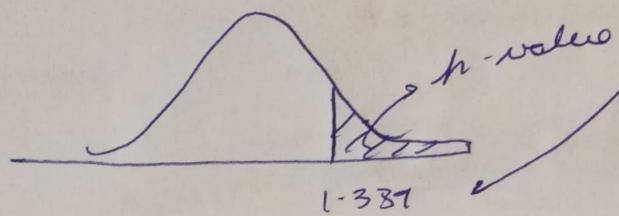
$$t = -1.31$$

Eg 25 dehydrated women. over $R_{20} = 1.3$

speed stat ... mean rev is differe from 0.95 sec.

$$\mu = 0.95 \quad s = 0.05 \quad t = \frac{0.18}{0.13/\sqrt{25}} = 1.389$$

25-1 DOF



Eg $\mu = 3.5$
 $\mu \neq 3.5$ both sides $\rightarrow 0.05$

$$\bar{x} = 3.3 \quad s = 0.851 \quad DF = 10 - 1 = 9$$

$$t = \frac{3.3 - 3.5}{0.851/\sqrt{10}} = -0.743$$

(-2.262 2.262) critical t

(-2.262 2.262)

If it falls in that region
 we fail to reject H_0 .

Module 1

EDA, structured data ... rectangular data,
DF, incases, NDS, estimation of deviation.
Mean } + weighted trimmed robust devt.
Med
excellency (deviations, variance, SD range, quartile)
percentile.

Data: → structured : pre-defined schema / specific concepts
rows, obs. units / variables.

unstructured :

freq, central tendency,
sample & population

Descriptive Statistics: Organizing & summarizing data

Inferential Statistics : formal methods for drawing some
good conclusions from data
(based on sample)

→ uses probability to determine how confident we
can be that conclusions are correct.

regression,
hypothesis

Mathematical Models: describe phenomena that occur in real
world.

① deterministic:

↳ when value is precisely determined
from another value.

↳ v, u, a

↳ high degree precision

Statistical Models: → used to predict life's more
uncertain situations

↳ based on idea that one value
affects another value.

↳ weather prediction (probability)

Probability: study randomness, mathematical tool.
→ chance of occurrence of event.

Population

→ collection.

SAMPLING

→ gain info about population

Categorical / qualitative

Numerical / quantitative

Pie Chart → 100% marking, data fits total.

Bar graph → ~~DO NOT~~ consider when data points are missing.

→ exact no.

→ size of category

MISSING DATA
→ sorted.

Sample:

↳ should have same characteristics as the population

⊗ Random Sampling

① Simple

② stratified

③ cluster

④ systematic

(pattern selection).

① Simple Random

- ↳ each has same chance of being selected.
- ↳ random no. generated $\frac{9436}{_}$ ②

② Stratified Sampling

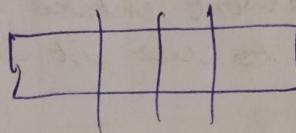
1	2	3	4
---	---	---	---

now perform random sampling
per each category.

(equal dist.)

May / may not.

③ Cluster



clusters

- select random clusters as a whole.
- may NOT contain equal no. of randomly chosen students from each class.

Non Random

① Convenience Sampling

- using results already available.
- sending fliers.

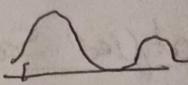
Sampling Errors

- ↳ sample size
- ↳ defective device \rightarrow non sampling errors.

Sampling bias: \rightarrow a sample is collected from population where ~~some~~ members of population do not have an equally likely chance of being chosen.



Problems

- ① sample - Biased
- ② self selected samples (internet etc.)
- ③ size sample
- ④ influential
- ⑤ refusal of participation → 
- ⑥ self centered
- ⑦ Misleading use of data - incorrect graphs etc.

NOTE:

- ① sampling types may cause variations.
- ② if it doesn't represent entire population
it is not useful to use it as sample

② Music store → convenience sample

200 ↗ 174 ←
→ 46 X NO NOT AD

Sampling Variability

→ samples taken by different parameters

→ large sample size → BETTER!!

Confidence intervals.

Scales

↳ nominal: 7 IN questions No 1 12

↳ ordinal: top 5; ordered hierarchical

↳ interval:

↳ Ratio:

CRV \rightarrow a RV whose outcomes are measured.

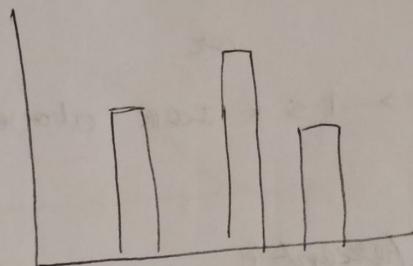
DISCRETE : result of counting

CONTINUOUS : result of measuring

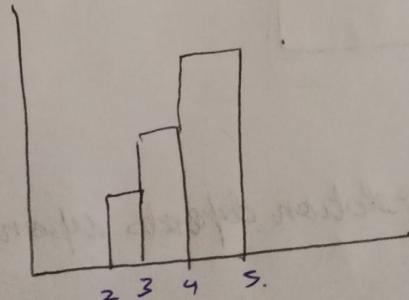
Sampling error - natural variation that results from selecting a sample to represent a larger population.

Bar-graphs

numerical data
 \rightarrow discrete
 \rightarrow frequencies on y



Histogram



\rightarrow continuous
 \rightarrow intervals
 \rightarrow tells the distribution
actually
 \rightarrow values within interval

Measures of location of data:

- ① quartiles
- ② percentile

1st Q 25th percentile

Median 2. 50th percentile.

Q₃ 75th percentile.

90^{th} percentile \Rightarrow 10% of the test scores are the same as or greater than your test score.

$$\text{IQR} \rightarrow Q_3 - Q_1$$

Outlier is determined by IQR.

$$x < 1.5 \times \text{IQR} \text{ below } Q_1$$

or

$$x > 1.5 \times \text{IQR} \text{ above } Q_3$$

Finding k^{th} percentile

$$i = \frac{k}{100} (n+1)$$

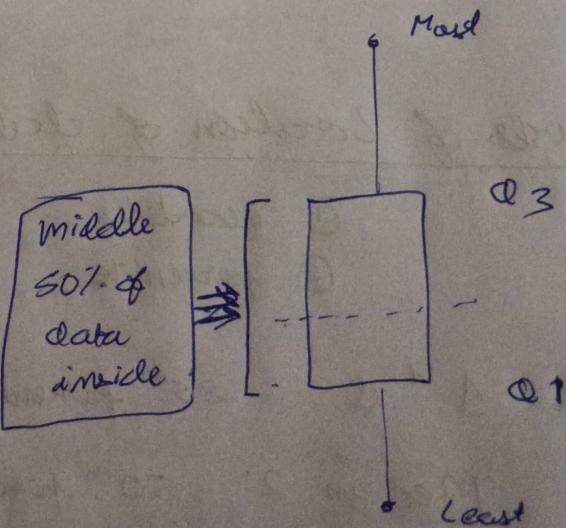
n : no. of data points

Percentile interpretation depends upon usage

Box Plots: ⑤

- median
- Q_1
- Q_3
- min
- max.

Middle
50% of
data inside



Measures of centre of data

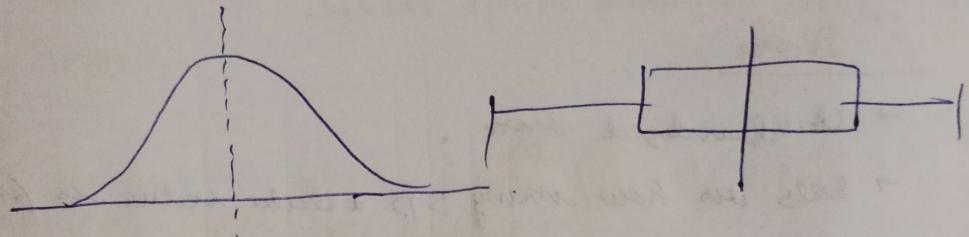
mean, median

median \rightarrow better when there are outliers

sample mean: \bar{x}

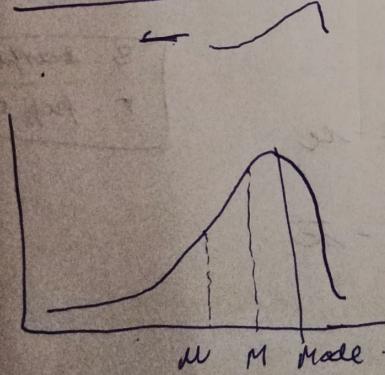
population mean: μ

Skewness



Mean / Median
(symmetric distribution)

LEFT



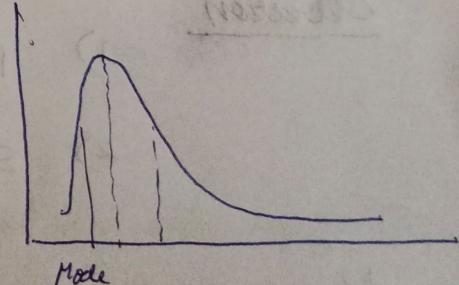
left tail being pull.

\rightarrow more higher values.

$\rightarrow \text{Mean} < \text{Median}$

$$\bar{x} < \text{Median} < \text{Mode}$$

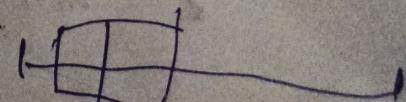
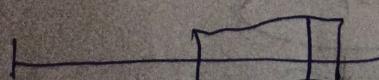
RIGHT



\rightarrow more lower values

$\rightarrow \text{Mean} > \text{Median}$

$$\bar{x} > \text{Median} > \text{Mode}$$



Measures of spread of data

→ variation in data

STANDARD deviation

- ① measure how far data values are from their mean.
 - ② determine if data pt is close/far from mean.
- It is ALWAYS +ve or zero.

Z-Score:

- standardized score.
- tells us how many S.D a data value is from the mean.

$$\bar{Z} = \frac{x - \mu_p}{\sigma_p}$$

x : particular score
 μ : population mean
 σ : population S.D

Deviation:

1) Population : $x_c - \mu$

S : sample S.D
 σ : pop SD

2) Sample : $x_c - \bar{x}$

Variance:

: average of squares of deviations.

$$S = \sqrt{\text{var}}$$

$$\text{var} = \frac{\sum (x_c - \bar{x})^2}{n-1}$$

$$\approx \frac{\sum (x_c - \bar{x})^2}{n-1}$$

sample

For population $\rightarrow n$

Bimodal

$$S = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N-1}} \rightarrow \text{Sample aka S.D}$$

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}} \rightarrow \text{population aka SD}$$

Sampling Variability of a statistic

↓ draw much statistic series from one sample to another.

measured by **[STANDARD ERROR]**

is the average S.D resulting from repeated sampling

$$\text{Standard error of mean} : \frac{\sigma}{\sqrt{n}}$$

σ : S.D population
n: size of sample

NOTE:

① mean affected the most by skewing.

$$\text{Mean absolute dev} = \frac{\sum |x_i - \bar{x}|}{N}$$

$$\text{M.A.D} = \text{Median}(|x_1 - m|, |x_2 - m|, \dots)$$

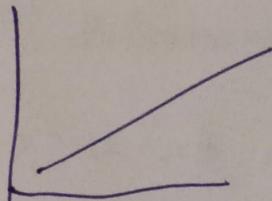
Probability

$$-1 \leq x \leq 1$$

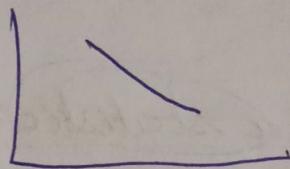
Expected value

$$E[X] = \sum_{x=1}^n p(x) \cdot x$$

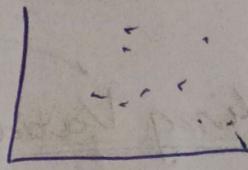
$$x = 0 \quad \text{no coree.}$$



+ve



-ve

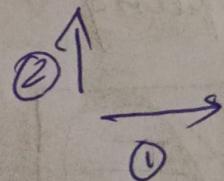


ONLY LINEAR

Hexagonal Binning:

explores 2 or more categorical variable as well.

Vidim Plot



Probability: certain we are of results....

Experiment:

operation carried under controlled conditions

↓ result

Outcome

sample Space: set of possible outcomes ...

Conditional

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Independent & MCE.

→ INDEPENDENT

Two events are independent if knowledge of occurring of one doesn't affect chance of other.

$$P(A \cap B) = P(A) \cdot P(B)$$

$$P(A|B) = P(A)$$

$$P(B|A) = P(B)$$

→ Mutually exclusive

either A or B.

$$P(A \cap B) = 0$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

A = first attempt success. 0.65

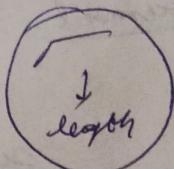
B = second attempt 0.65

$$P(CB/A) = 0.90.$$

$$A \cap B = ?$$

$$0.90 \times 0.65.$$

↓
both goals.



t_{dist}
 $N > 30$
normal
 $N < 30$
 t_{dist}

either $P(A \cup B)$

any door is $\frac{1}{3}$

Eg: 3 doors.

1 out \rightarrow

caught

$$\frac{1}{3}$$

not caught.

$$\frac{2}{3}$$

$$\frac{5}{15}$$

2 out \rightarrow

$$\frac{1}{4}$$

$$\frac{3}{4}$$

$$\frac{3}{12}$$

3 out \rightarrow

$$\frac{1}{2}$$

$$\frac{1}{2}$$

$$\frac{1}{6}$$

$$\frac{1}{60}$$

$$\frac{81}{60}$$

$$1$$

$$\frac{\frac{1}{15} + \frac{1}{12}}{29/60}$$

$$\frac{9}{60} = \frac{3}{20}$$

Eg: $P(\text{bitter}) = 0.2$

$\text{no bitter} = 0.8$

$P(\text{bark} | \text{bitter}) = 3\%$

$$P(A \cap B) = P(A \cap B) / P(B)$$

$P(\text{bark} | \text{no bitter}) = 4\%$

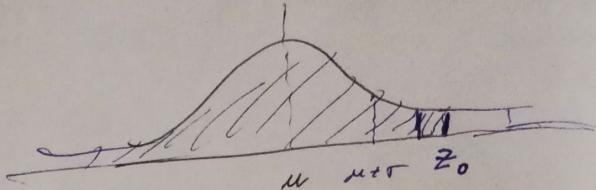
$P(\text{bitter} | \text{bark}) = ?$

$$P(\text{bark} \cap \text{bitter}) = 3\% \cdot 0.2 = P(\text{bitter} | \text{bark}) \cdot P(\text{bark})$$

$$P(\text{bark} \cap \text{no bitter}) = 4\% \cdot 0.8 = P(\text{no bitter} | \text{bark}) \cdot P(\text{bark})$$

Normal distribution

$$X \sim N(\mu, \sigma^2)$$



→ depends only on mean & ~~median~~ S.D

→ area = 1

Standard Normal distribution

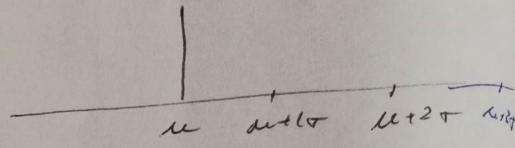
$$\mu = 0$$

$$\text{S.D} = 1$$

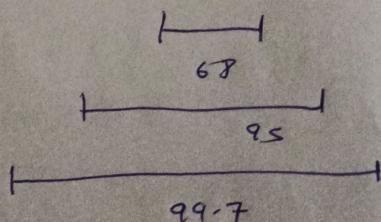
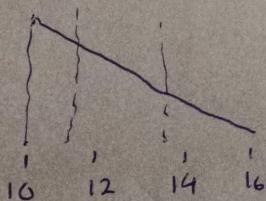
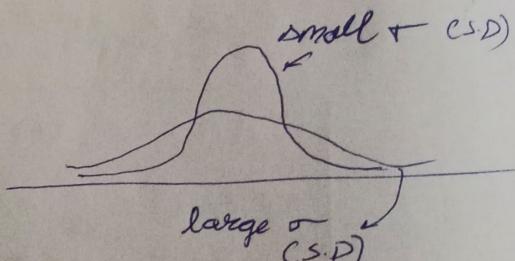
$$N(0, 1)$$

Z-score: indicates no. of S.D a score falls above/below the mean.

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



$$P(Z \leq z_0) = \text{area until } z_0$$



$$Z_{11} = \frac{x - \mu}{\sigma}$$

$$= \frac{11 - 10}{2} = 0.5$$

$$Z_{13.6} = \frac{3.6}{2} = 1.8$$

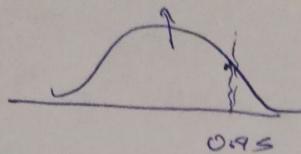
$$\Phi = \alpha_{z=1.8} - \alpha_{z=0.5}$$

Eg:

Top 5%.

N(500, 100)

$x = ?$



$$\text{area} = 0.95$$

$$\Rightarrow z = 1.645$$

$$z = \frac{x - \mu}{\sigma}$$

$$1.645 \times 100 + 500 = x$$

$$x = 664$$

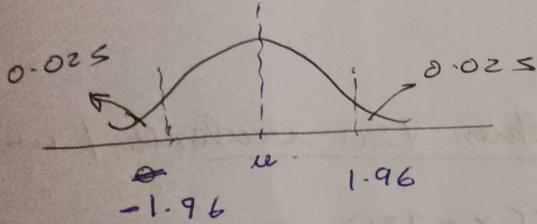
Sampling Distribution

Standard dev of
mean

$$\sigma_{\bar{x}} = \frac{\sigma_{\text{Pop.}}}{\sqrt{m}}$$

σ : population

m : sample size



if $z \leq 1.96$
within what lim
sample average $P = 0.95$

$$X = \mu - 1.96 (\text{S.D. mean})$$

$$\mu + 1.96 (\text{S.D. mean})$$

In Normal distribution $\text{Mean} = \text{Median}$

In continuous distribution $P(X < 1) \approx P(X \leq 1)$

Very Very IMP

Ex: 6000 weekly
10% defective

random sample $n = 100$

X : no. of defective in a

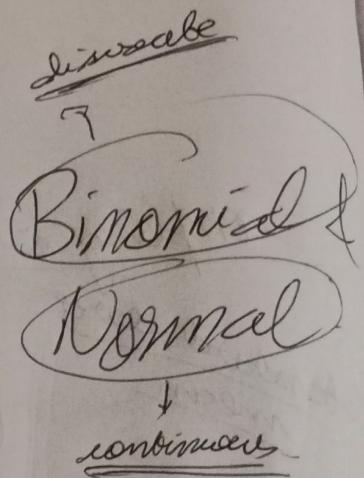
$X = ?$

$$n = 100$$

$$p = 0.1$$

$$q = 0.9$$

$$\boxed{\begin{aligned} \mu &= np \\ \sigma &= \sqrt{npq} \end{aligned}}$$



Ex: 75% \rightarrow favour.

find probability more than 780 students / 1000

will be in favour. $P(X = 781) + P(X = 782) \dots$

$$\mu = 0.75$$

$$\mu = np = 750.$$

$$\sigma = 0.25$$

$$\sigma = \sqrt{npq} = 13.7.$$

$$n = 1000.$$

continuous

$$N(750, 13.7^2)$$

$$P(780.5 < X < 1000.5)$$

$$\approx P\left(\frac{780.5 - 750}{13.7} \leq Z \leq \frac{1000.5 - 750}{13.7}\right) \quad \text{Z-score}$$

Eg:

defective 10%.

$$n = 200$$

P of defective 24 & 30

$$np = 20$$

$$q = .9$$

$$n = 200$$

$$\mu = np = 20$$

$$\sigma = \sqrt{npq} = 4.24$$

$$P(25 \leq X \leq 29)$$

$$\approx P(24.5 \leq X \leq 29.5)$$

$$P\left(\frac{|20 - 24.5|}{4.24} \leq Z \leq \frac{|29.5 - 24.5|}{4.24}\right)$$

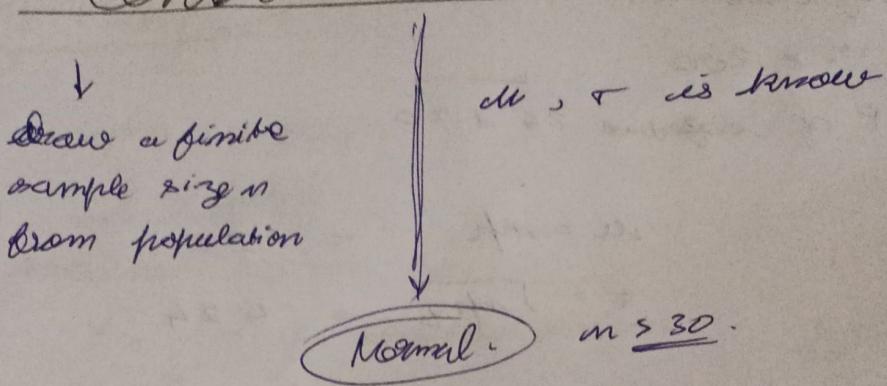
$$Z_1 = -1.18$$

$$Z_2 = +1.18$$

$$Z_1 = -0.27$$

Central Limit Theorem

D)



μ_n : mean of population

σ_n : S.D of population then

$$X \sim N\left(\mu_n, \frac{\sigma_n}{\sqrt{n}}\right)$$

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

sample size

Sample means form their own

$$n \geq 30$$

Ex: $\mu_{\text{pop}} = 90$

$$n = 25$$

$$\sigma_{\text{pop}} = 15.$$

a)

for Normal distribution

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

for sample d.

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

Ex: $\mu = 74$ } normal distribution.
 $\sigma = 6.8.$

a) selected @ random $\bar{x} < 65$.

$$P\left(Z < \frac{65 - 74}{6.8}\right) \leftarrow \text{area}$$

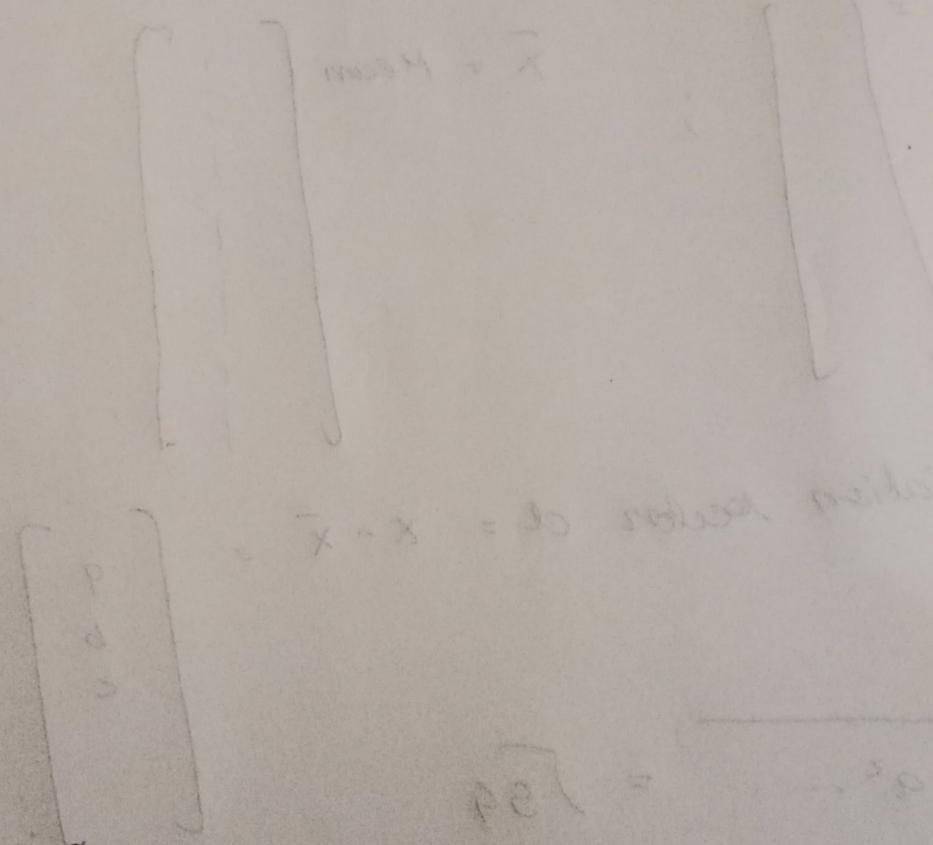
D) A sample of 50 selected @ exam $P(Z > 7.5)$

$$n = 50$$

$$P\left(Z > \frac{75 - 74}{\frac{6.8}{\sqrt{50}}}\right)$$

mean exam
score

MEAN \rightarrow CLT.



CIE

Eg: Deviation data ... length of deviation vector ...

$$x = \underbrace{[\quad]}_{\text{some 10 pts.}}$$

① find mean

$$\textcircled{2} \quad x = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}; \quad \bar{x} = \text{Mean} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$\textcircled{3} \quad \text{deviation vector } d = x - \bar{x} = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

$$\|d\| = \sqrt{a^2 + \dots} = \sqrt{S_4}$$

Histogram \rightarrow tells the distribution ... can be plotted in boxplot.

\rightarrow continuous data. (median, quart ...)

Bar graph \rightarrow categorical

Eg: Disjoint are M-E

Hypothesis Testing Qs.

Eg:

$$\mu = 1.2$$

(not) \sim

P(0)

$$\mu \neq 1.2$$

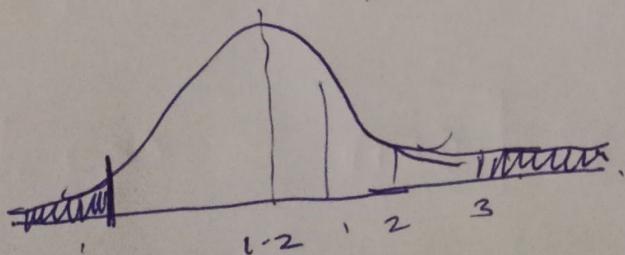
24.5

t, P
 $p < 0.05$
 $\alpha = 0.05$
 $m = 2.50$
 26.05 ± 0.50

the mean of 100 injected is 1.05 ~~and~~ S.D = 0.5

Remember:

$$\sigma_{\bar{x}} = \frac{s}{\sqrt{n}} = 0.05$$



$$z = \frac{1.2 - 1.05}{0.05}$$

$z = 3-$

Random Variable

DISCRETE : counting

continuous : ht., temp. measuring

Probability Distribution Function (PDF)

Discrete

$$\textcircled{1} \quad 0 \leq P(x) \leq 1$$

$$\textcircled{2} \quad \sum P(x) = 1$$

} conditions

Expected Value: $E(x)$ Mean.

$$E(x) = \sum x P(x)$$

long term average

Binomial distribution

3 characteristics.

① fixed no. of trials. (n)

$$n \times k \times (1-p)^{n-k}$$

② ONLY two possible outcomes

p : success

q : failure

$$p + q = 1$$

③ the trials are independent

↳ should be consistent
for all trials.

If $n=1$

↳ Bernoulli trial

+②

+③

Mean, $\mu = np$

Variance $\sigma^2 = npq$

S.D $\sigma = \sqrt{npq}$

$$P(X=x) = p^x (1-p)^{n-x}$$

$x = 0, 1$

Geometric Distribution "first"

$$1 - \frac{7c_2}{10} = \frac{7c_2}{270}$$

- ① repeating independent Bernoulli trials until success is obtained.

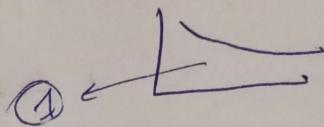
- ② atleast 1 trial

- ③ p, q

$$p(x) = (1-p)^{x-1} p^1$$

$$\text{Var} = E(x^2) - [E(x)]^2$$

$\hookrightarrow p = \text{success}$.



Poisson

- ① no. of occurrence of an event in a given unit of time, distance / area / volume.

$$P(X=x) = \frac{\lambda^x e^{-\lambda}}{x!}$$

λ : mean time.

$$P(X \leq s) = P(X=0) + P(X=1) + \dots + P(X=s) \\ = q^0 p^0 + q^1 p^1 + q^2 p^2 + \dots$$

e.g. 2% defect \rightarrow 100 bypes

- a) $P(X < 7)$ within first bypes \rightarrow defective piece.
- b) $P(X > 12)$
- c) $P(X \geq 8)$
- d) $P(15 \leq X \leq 30)$
- e) $P(20 < X \leq 45)$

a) $P(X < 7) \rightarrow 1 - q^7$

$$P = 1 - (0.98)^7$$

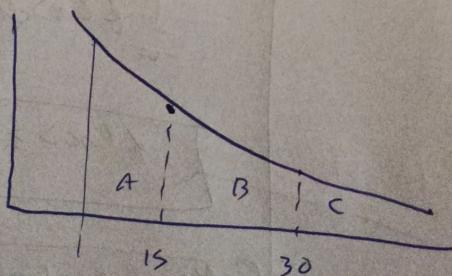
$$P(X \geq x) = 1 - q^x$$

b) $X \leq 12 \quad 1 - (0.98)^{12}$

$$\text{BUT } x > 12 \Rightarrow 1 - (1 - 0.98)^{x-12}$$

c) $X \leq 17 \quad 1 - (0.98)^{17}$

$$(0.98)^{17} =$$



$$P(X \leq 15) + P(15 \leq X \leq 30) + P(X \geq 31) = 1$$

$$P(X \leq 1) + P(X \geq 1) = 1$$

Eg.

12 customers / day

- a) what is probability exactly 8 customers in 1 day
 \rightarrow freq. given Poisson !!

$$\mu = 12$$

$$P(X=x_c) = \frac{\mu^x e^{-\mu}}{x!}$$

$$\begin{aligned} & \left(\frac{1}{2}\right)^n \times \left(\frac{1}{2}\right)^m \times \frac{2^n}{2!} \\ & \left(\frac{1}{2}\right)^8 \times \left(\frac{1}{2}\right)^4 \times \frac{2^8}{2!} \\ & \frac{1}{2^{12}} \times \frac{1}{2^4} \end{aligned}$$

$$P(X=8) = \frac{12^8 e^{-12}}{8!}$$

Eg.

average \rightarrow 7 bent in 2 hrs.

- a) exactly 9 bent in 2 hrs.

$$\mu = 7$$

$$x = 9$$

$$P(X=9) = \frac{7^9 e^{-7}}{9!}$$

- b) exactly 24 bent in 8 hrs.

$$\mu = 7 \times \frac{8}{2} = 28$$

$$P(X=24) = \frac{28^{24} e^{-28}}{24!}$$

PMF

$$\sum P(x) = 1$$

$$P(x) \geq 0$$

ALWAYS

PDF

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

$$f(x) \geq 0$$

$$\text{Mean doftime} = \mathbb{E}[x] = \int x f(x) dx$$

Valid PDF

$$\int_6^8 f(x) = 1$$

	mean	s.d.
+3	+3	-
x2	x2	x2
x-1	x-1	-

E.g. given day forgetting = 3

Bayes

a) $P(X=5) = \frac{3^5 e^{-3}}{5!}$

b) $P(X=5 | X \geq 4) = \frac{P(X=5)}{P(X=4) + P(X=5)}$

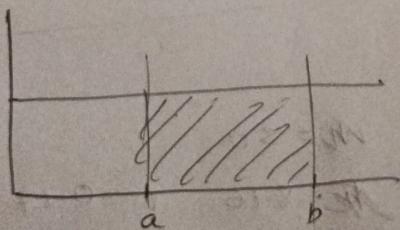
Covariance

CRV

↳ uniform distribution

↳ exponential distribution

Uniform distribution



$$f(x) = \frac{1}{b-a}$$

$$P(X = \frac{a+b}{2}) = 0.$$

$$\text{Mean} = \frac{a+b}{2}$$

$$\text{var} = \frac{(b-a)^2}{12}$$

You flip a coin ... $\mu(H) = 50\%$ is known but think it is less ... Type of test?

$$n = 50 \quad \bar{x} = 0.46$$

$$\alpha = 0.01$$

$$\mu = 0.5$$

Null

$$\bar{x} < 0.5$$

Alternative

Verify - given

2. S.D.

single tailed

$$Z = \frac{\bar{x} - \mu}{\sqrt{\frac{\mu(1-\mu)}{n}}}$$

$$\frac{\bar{x} - \mu}{\sqrt{\frac{\mu(1-\mu)}{n}}}$$

$$\sqrt{\frac{\mu(1-\mu)}{n}}$$

$$= \frac{0.46 - \mu}{\sqrt{0.5 \times 0.5 \times 50}}$$

IMP.

$$\text{Var}(X - aY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)$$

$$\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$$

ab

$$H_0: \mu = 8.21$$

$$n = 50$$

$$H_1: \mu < 8.21$$

$$\bar{x} = 8.16 \quad t = 0.17$$

$$\alpha = 0.01$$

$$Z = \frac{8.16 - 8.21}{\sqrt{0.17 / 50}} = -2.08$$

$$Z_{-0.01} = -2.33$$

FTR H₀

