



Non-Parametric tests :

(Parametric Tests) -
Conventional & fair

Test-statistic

$$- \frac{\bar{X} - c}{\sigma/\sqrt{n}} ; \quad \frac{\bar{X} - c}{S/\sqrt{n}} ; \dots$$

- to perform the test assumption of

Sample size
being large

Normal distribution or by

assuming C.L.T. holds

Non-parametric test :-

No need to
Device a
test
for
normality

- Distribution free tests
 - { - less powerful than parametric tests
- Power of test :- next statistics course
- Very helpful for a start



Example :- - Vaccine has been developed

Q:- Is the vaccine effective or not?
(for a disease)

A:- Designed an experiment

- choose n individuals from a population

- gave vaccine to n_1 of them
 placebo to $n_2 := n - n_1$

	Affected	Not affected	Total
vaccine	X_{11}	X_{12}	n_1
placebo	X_{21}	X_{22}	n_2

\hat{Q} : Based on table can we answer 150 question Q ?

Example 2 :- [Tea-Tasting Example]

A-claimed :- Can tell from tasting a cup of tea whether :-

- English preparation
- milk first and Tea next
 - Tea first and milk next

Designed an Experiment :-

- Prepared 8 cups of tea
-

- Person A tasted each one of them and gave opinion.

	Tea	Milk	Total
Tea first	3	1	4
Milk first	1	3	4

Approach 1 (Parametric) : χ^2 - test for independence.

Approach 2 :-

Experiment:	Tea	Milk	Total
Tea first	$3 \equiv X_{11}$	1	4
Milk first	1	3	4

- Row totals are fixed by the Experimenter

- Column totals are decided by Person A

H_0 :- Person A has No ability
 (i.e. Person A calls "tea first"
 by choosing a random sample
 from S)

Under H_0 :-

$$P(X_{11} = 3) = P(\text{choose 4 cups from 8 cups } \in \text{5 of them are Tea})$$

$$(\text{Hypergeometric}) = \frac{\binom{4}{3} \binom{4}{1}}{\binom{8}{4}} = 0.229$$

Test p-value :- $P(X_{11} \geq t_0)$ observed

Significance level := $\alpha \leq 0.05$

Here $t_0 = 3$:- Assume H_0 is true

$$P(X_{11} \geq 3) = P(X_{11} = 3)$$

$$+ P(X_{11} = 4)$$

Under H_0 :- $X_{11} \sim \text{Hypergeometric}(N=8, G=4, n=4)$

$$P(X_{11} \geq 3) = 0.229 + \frac{\binom{4}{4} \binom{4}{0}}{\binom{8}{4}}$$

$$= 0.229 + 0.014$$

$$= 0.243$$

∴ As $P(X_{11} \geq 3) = 0.243 > d=0.05$

∴ the null cannot be rejected.

- there is not enough evidence to reject

the null hypothesis (that person

A was purely guessing

Sign test & Signed Rank tests

Model : X_1, \dots, X_n are from

a random sample

$x_i = \theta + \epsilon_i$ Errors
 independent with
 p.d.f + C. Symmetric
 around 0.

$$H_0: \theta = 0 \quad H_A: \theta > 0$$

Test statistic

$$S = \sum_{i=1}^n sgn(x_i)$$

where $sgn(t) = \begin{cases} -1 & t < 0 \\ 0 & t = 0 \\ 1 & t > 0 \end{cases}$

Sign test

- look at positive observations

$$S^+ = \#\{i : x_i > 0\}$$

- ignore 0's and sample is reduced.

- If we observe -1 and 1 then

X_i - can be thought as Bernoulli(.)

$\Rightarrow S^+ \sim \text{Binomial distribution}$
 $(n, \frac{1}{2}) = \text{Null distribution}$

Test :- Compare this distributions with
the observed statistic

Example (Sign test) :-

12 people are chosen

10 prefers shorts over full pants $(X_i > 0)$

1 prefers full pants are shorts

1 no preference $(X_i = 0)$ $(X_i < 0)$

- strong preference for shorts ?

How likely is such a result true if

H_0 : There is no preference over shorts or full pants

is true ?

$$\delta^+ = \{ i : X_i > 0 \sim \text{Bin}(n=11, p=\frac{1}{2}) \}$$

$$g^+ \equiv \text{observed} = 10$$

P-value :- $\mathbb{P}(\delta^+ \geq 10) \approx 0.0059$

α - level of significance :- 0.05

$$\Rightarrow \text{Since } 0.0059 << 0.05$$

H_0 : No preference H_A : Shows overfull

we reject null hypothesis.

Signed Rank - Wilcoxon Test

(Compare with t-test)

$$H_0: \mu = 0, \quad H_A: \mu > 0$$

Not distribution free

$$\text{Test-statistic} \cdot t_0 := \frac{\bar{X}}{S/\sqrt{n}} \sim t_{n-1}$$

$$\text{P-value} : \mathbb{P}(t_{n-1} > t_0)$$

x_1, x_2, \dots, x_n - sample $\{x_i = \theta + \epsilon_i \text{ (old model)}$

$$W = \sum_{i=1}^n \text{sgn}(x_i) \text{Rank}(|x_i|)$$

Test-statistic

$$W^+ = \sum_{j=1}^n \text{Rank}(|x_j|) \mathbb{1}(x_j > 0)$$

W^+ symmetric
around $n(\bar{x} + 1)$

$$\stackrel{\text{Ex}}{=} \frac{1}{2} \left(W^+ + n(\bar{x} + 1) \right)$$

- There is no closed form for distribution of W^+

$H_0: \theta = 0$ versus $H_A: \theta > 0$

p-value: $P(W^+ > t_0)$

t_0 - observed test statistic

Recall - clean up

X_1, X_2, \dots, X_n are i.i.d Bernoulli(n, p)

$$\Rightarrow S_n \sim \text{Binomial}(n, p)$$

$$[\text{CLT}] \quad \frac{S_n - np}{\sqrt{np(1-p)}} \xrightarrow{d} N(0,1)$$

Approximation - Application : $a < b \in \mathbb{R}$

$$P(a < S_n \leq b) = P\left(\frac{a - np}{\sqrt{np(1-p)}} < \frac{S_n - np}{\sqrt{np(1-p)}} \leq \frac{b - np}{\sqrt{np(1-p)}}\right)$$

$$\approx P\left(\frac{a - np}{\sqrt{np(1-p)}} < Z \leq \frac{b - np}{\sqrt{np(1-p)}}\right)$$

when $Z \sim N(0,1)$

- When approximation is valid

$$- np > \frac{1}{2} \quad n(1-p) > \frac{1}{2}$$

- $n \ggg$ large

To make approximation precise / more accurate

- Continuity correction

Example :- $n=20, P=\frac{1}{2} \quad S_n \sim \text{Binomial}(20, \frac{1}{2})$

$$\mathbb{P}(8 \leq S_n \leq 10) = \mathbb{P}\left(\frac{8 - np}{\sqrt{np(1-p)}} \leq \frac{S_n - np}{\sqrt{np(1-p)}} \leq \frac{10 - np}{\sqrt{np(1-p)}}\right)$$

Central
limit
Theorem
approximation

$$\approx \mathbb{P}\left(\frac{8 - np}{\sqrt{np(1-p)}} < z \leq \frac{10 - np}{\sqrt{np(1-p)}}\right)$$

$$\approx \mathbb{P}\left(\frac{8 - 10}{\sqrt{5}} < z \leq \frac{10 - 10}{\sqrt{5}}\right)$$

$$= \Phi(0) - \Phi\left(-\frac{2}{\sqrt{5}}\right)$$

$$= \frac{1}{2} - \dots = 0.3145$$

Table or R

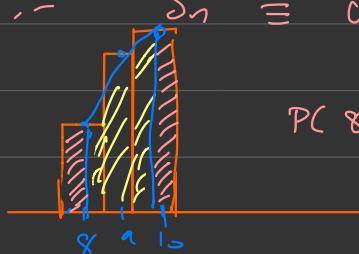
$$\mathbb{P}(8 \leq S_n \leq 10) = \sum_{k=8}^{10} \binom{20}{k} \left(\frac{1}{2}\right)^{10}$$

is not great

Direct
computation

$$= 0.4565$$

Cause for Bad Approximation :- S_n is discrete and z is continuous



$\mathbb{P}(8 \leq S_n \leq 10)$

Approximation - Application

with Continuity Correction : $a < b \in \mathbb{R}$

$$\mathbb{P}(a < S_n \leq b) = \mathbb{P}\left(a \frac{-np}{\sqrt{np(1-p)}} < \frac{S_n - np}{\sqrt{np(1-p)}} \leq \frac{b - np}{\sqrt{np(1-p)}}\right)$$

$$\approx \mathbb{P}\left(\frac{a - \frac{1}{2} - np}{\sqrt{np(1-p)}} < Z \leq \frac{b + \frac{1}{2} - np}{\sqrt{np(1-p)}}\right)$$

when $Z \sim N(0,1)$

$$\mathbb{P}(8 \leq S_n \leq 10) = \mathbb{P}\left(\frac{8 - np}{\sqrt{np(1-p)}} \leq \frac{S_n - np}{\sqrt{np(1-p)}} \leq \frac{10 - np}{\sqrt{np(1-p)}}\right)$$

Central
limit
Theorem
approximation

$$\approx \mathbb{P}\left(\frac{7.5 - np}{\sqrt{np(1-p)}} < Z \leq \frac{10.5 - np}{\sqrt{np(1-p)}}\right)$$

$$\approx \mathbb{P}\left(\frac{7.5 - 10}{\sqrt{5}} < Z \leq \frac{10.5 - 10}{\sqrt{5}}\right)$$

with Continuity

Correction

$$\Phi\left(\frac{0.5}{\sqrt{5}}\right) - \Phi\left(-\frac{2.5}{\sqrt{5}}\right) \\ = 0.4567$$

This is indeed a better approximation!

$X_i : i=1, \dots, n$ are discrete

$$S_n = X_1 + X_2 + \dots + X_n$$

Continuity correction

Continuity correction

$$P(a \leq S_n \leq b) = P(a - \frac{1}{2} \leq S_n \leq b + \frac{1}{2})$$

Histogram
or if $a, b \in \mathbb{N}$

CDF ... approximation ...

Boot strap - Jack-knife :- [R sampling techniques]

Efron -

- Parameter of interest $\theta \in \mathbb{P}$
- Estimate the parameter $\hat{\theta}$ Method of moments
Maximum likelihood estimate
- Confidence interval j - Hypothesis test
 χ^2 -test, ... : parametric
(Normal approximation)

Setting:- Normal approximation does not apply

Random Sample. $X_1, \dots, X_n \sim F$ ($F \equiv$ distribution of population)

θ - is parameter of interest

Compute :- $\hat{\theta} := g(X_1, X_2, \dots, X_n)$ is an estimate for θ .

$\hat{\theta} \sim G_n$ (some sampling distribution)

G_n :- Unknown if normality holds $\Rightarrow h_n \sim N(\cdot, \cdot)$

| (Ideal) Find as many datasets (Samples)
from G_n (B)

Wish
more samples
of G_n .

then $\hat{\theta}_1, \dots, \hat{\theta}_B \sim G_n$.

Once we have B realisations from G_n

→ estimate many characteristics of G_n

Bootstrap :- - Our current data set is

$$x_1, x_2, \dots, x_n$$

Make it precise

↓

- Resample repeatedly from data set with replacement.

(Hope ...) - Approximate sample from G_n . & using this we could estimate characteristics of the distribution of G_n

Non-parametric Bootstrap :-

- $\{x_1, \dots, x_n\} \equiv S$ - sample from population.
 $\theta := \phi(x_1, \dots, x_n)$

- Resample from S WITH replacement a sample of size n .

$$x_{11}^*, \dots, x_{n1}^* \Rightarrow \hat{\theta}_1^*$$

Repeat
 B
time

$$x_{12}^*, \dots, x_{n2}^* \Rightarrow \hat{\theta}_2^*$$

$$\vdots \quad \vdots \quad \vdots$$
$$x_{1B}^*, \dots, x_{nB}^* \Rightarrow \hat{\theta}_B^*$$

B
Estimation

$$(\text{Ex.}) \quad \hat{\theta}_1^*, \dots, \hat{\theta}_B^* \approx g_n$$

• $\text{Var } \hat{\theta} = \frac{1}{B-1} \sum_{i=1}^B (\hat{\theta}_i^* - \bar{\hat{\theta}}_B^*)^2$

Boot strap estimate Sample mean

• Boot strap confidence interval

$$(\text{order}) \quad \hat{\theta}_{(1)}^* < \dots < \hat{\theta}_{(k)}^* < \dots < \hat{\theta}_{(B)}^*$$

$$L = \hat{\theta}_{\lfloor \frac{\alpha}{2} B \rfloor} \quad U = \hat{\theta}_{\lfloor (1-\frac{\alpha}{2}) B \rfloor}$$

$\Rightarrow (L, U)$ is $100(1-\alpha)\%$ Boot strap
confidence interval for θ

Parametric Boot strap

$$x_1, \dots, x_n \sim F(\theta)$$

↙ unknown
 ↘ known distribution

$$S = \{x_1, \dots, x_n\}$$

set $\hat{\theta} = g(x_1, x_2, \dots, x_n)$

- Resample from $F(\hat{\theta})$ WITH replacement - a sample of size n .

$$x_{11}^*, \dots, x_{n1}^* \sim F(\hat{\theta}) \implies \hat{\theta}_1^*$$

Repeat

B

$$x_{12}^*, \dots, x_{n2}^* \sim F(\hat{\theta}) \implies \hat{\theta}_2^*$$

time

$$\vdots$$

$$x_{1B}^*, \dots, x_{nB}^* \sim F(\hat{\theta}) \implies \hat{\theta}_B^*$$

B
Estimation

$$\frac{1}{B} \sum_{i=1}^B \hat{\theta}_i^* := \bar{\hat{\theta}}^*$$

- Estimate for $\hat{\theta}$

- Same procedure for confidence interval

Example : Population with distribution $\sim \text{Gamma}(\alpha_1)$

Sample :- x_1, \dots, x_n from population

$$\hat{\theta} := \bar{X} = \frac{\sum_{i=1}^n x_i}{n} \text{ is an estimate for } \alpha.$$

Bootstrap

Non-parametric :- Resample from $S = \{x_1, \dots, x_n\}$

Parametric :-

$$x_{11}^*, \dots, x_{n_1}^* \sim \text{Gamma}(\bar{x}, 1) \rightarrow \bar{x}_1^*$$

$$x_{12}^*, \dots, x_{n_2}^* \sim \text{Gamma}(\bar{x}, 1) \rightarrow \bar{x}_2^*$$

$$\vdots \quad \vdots$$

$$x_{1B}^*, \dots, x_{n_B}^* \sim \text{Gamma}(\bar{x}, 1) \rightarrow \bar{x}_B^*$$

$\frac{\alpha}{2}, \frac{1-\alpha}{2}$ quantiles of \bar{x}_i^* \Rightarrow confidence intervals

Central limit Theorem [check : mean & variance of $\text{Gamma}(a, 1)$]

$$\sqrt{n} \left(\bar{x} - a \right) \xrightarrow{d} N(0, 1)$$

confidence interval :

$$\left(\bar{x} - z_{\frac{1-\alpha}{2}} \sqrt{\frac{\bar{x}}{n}}, \bar{x} + z_{\frac{1-\alpha}{2}} \sqrt{\frac{\bar{x}}{n}} \right)$$

4 Compositions :-

- non parametric distributions

R
code

- Parametric distribution

- Normal distribution approximation

- True sampling distributions

