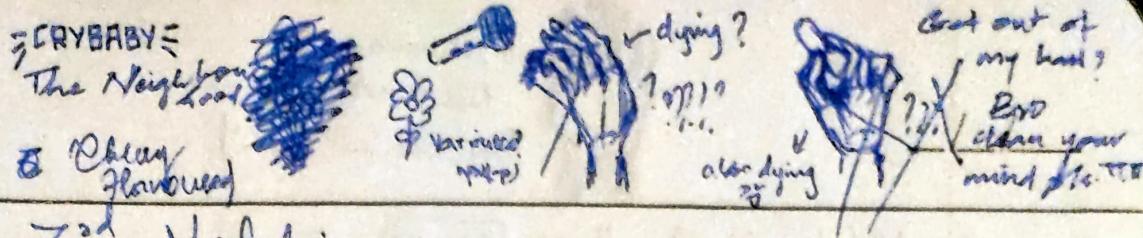


13 Feb.



3rd Module:

* Statistic Preparation:

p - represents the population's proportion.

\hat{p} - represents the sample's proportion.

$$\hat{p} = \frac{x}{n} \rightarrow \# \text{ individuals in the sample with the characteristic.}$$

$\rightarrow \# \text{ individuals in the sample.}$

$$\hat{p} \rightarrow [0, 1].$$

→ binomial.

$$\text{Mean of } \hat{p} = E[\hat{p}] = E\left[\frac{x}{n}\right] = \frac{1}{n} E[x] = \frac{1}{n} \times np = p.$$

$$\begin{aligned} \text{Variance of } \hat{p} = \text{Var}\left(\frac{x}{n}\right) &= \frac{1}{n^2} \text{Var}(x) = \frac{1}{n^2} \times np(1-p) \\ &= \frac{1}{n} p(1-p) = \frac{p(1-p)}{n}. \end{aligned}$$

$$\text{Standard deviation of } \hat{p} = \sqrt{\frac{p(1-p)}{n}}.$$

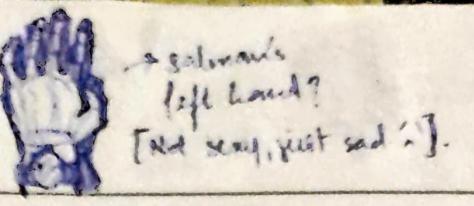
Central Limit Theorem: Condition $np \geq 15$. OR $n(1-p) \geq 15$.

$\uparrow n \text{ & } p$. Even when n is small, if $p = 0.5$, it gives
a \rightarrow a considerably normal distribution.

With n increasing, it draws closer to a normal distribution with smooth.

When n is small, $p = 0.09$, it gives \rightarrow a left skewed graph. Only with increase in size n can make it normal.

Similarly, when $p = 0.96$, it gives \rightarrow a right skewed graph. With increase in n , we get more normal dist.



$$N(\mu, \sigma) \Rightarrow N(p, \sqrt{p(1-p)/n})$$

Note →

Normal distribution ~~only~~ \Rightarrow for p too low or too high ^{only}, when the value of n is large.

- Qn. In a certain liberal precinct(n) across Town, 81% of the voters are registered democrat. What is the probability that, in a random sample of 100 voters from this precinct, no more than 80% of the voters would be registered democrat.

$$\rightarrow n = 100. \quad \hat{p} = \frac{80}{100} = 0.8. \quad p = 81\% = 0.81.$$

$$\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}} = \sqrt{\frac{0.81(0.19)}{100}} = 0.039.$$

$$z = \frac{\hat{p} - p}{\sigma_{\hat{p}}} = \frac{0.8 - 0.81}{0.039} = \frac{-0.01}{0.039} = -0.256$$

0.05 | 0.06
-0.2 0.901 0.397

$$\therefore P(\hat{p} \leq 0.8) = P(z \leq -0.256) = 0.3999.$$

Confidence Interval for p :

$$\hat{p} \pm z_{\alpha/2} \cdot \sigma_{\hat{p}}$$

$$\text{where standard error } \sigma_{\hat{p}} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}.$$

Qn. A study investigated the proportion of male birth in Liverpool, England. In a sample of 5045 births to non-smoking parents, there were 2685 males and 2360 females born. Construct a confidence interval for p, i.e., pop proportion of male birth.

$$\hat{p} = \frac{2685}{5045}$$

$$\hat{p} = \frac{2685}{2685 + 2360} = \frac{2685}{5045} = 0.5322$$

$$SE(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \sqrt{\frac{0.5322(0.4678)}{5045}} = \sqrt{0.00004935} \\ = 0.007025$$

$$\therefore CI = \hat{p} \pm z_{\alpha/2} \cdot SE(\hat{p}) = 0.5322 \pm 1.96 \cdot 0.007025 \\ = 0.5322 \pm 0.013769 \\ \rightarrow \underline{(0.518431, 0.545969)}$$

Qn. A recent survey of 1000 married men revealed that 561 of them have been unfaithful atleast once. Form a 95% confidence interval to estimate the true proportion of married men who are unfaithful.

$$\hat{p} = \frac{561}{1000} = 0.56$$

$$SE(\hat{p}) = \sqrt{\frac{0.56(0.44)}{1000}} = 0.0157$$

$$CI = \hat{p} \pm z_{\alpha/2} \cdot SE(\hat{p}) = 0.56 \pm 1.96 \times 0.0157 \\ = 0.56 \pm 0.03072 \Rightarrow \underline{(0.529278, 0.59072)}$$

$$\therefore CI = \underline{(0.529278, 0.59072)}$$

24 Feb.

* Sampling Variability for CI and Required size!

Suppose we are about to draw a sample and also wish to estimate the population proportion p . We may wish to estimate p to within an amount ϵ with 95% confidence.

How large of a sample size is required?

$$\rightarrow 1.96 \times \sqrt{\frac{p(1-p)}{n}} \leq \epsilon$$

max when: $\underline{p = 0.5}$.

Q. * Note from others:

* Hypothesis Testing:

$H_0: p = p_0$. \rightarrow Hypothesized value of p .

$H_1: p > p_0$ or $p < p_0$ or $p \neq p_0$.

$$Z = \frac{\hat{p} - p_0}{\text{SE}_{\text{p}}(\hat{p})}$$

where $\text{SE}_{\text{p}}(\hat{p}) = \sqrt{\frac{p_0(1-p_0)}{n}}$

Depends on H_1 :

When $H_1: p < p_0$. Left tailed.

$H_1: p > p_0$. Right tailed.

$H_1: p \neq p_0$. Two tailed.

If p value $\leq \alpha$, reject H_0 .

On ~~people~~ birth birth
In a random sample of 200 parts from a manufacturing process, 18 had major defects on male birth.

$$\hat{p} = \frac{18}{200} = 0.09.$$

Test the null hypothesis that the true proportion of male births is 0.5, against the alternative hypothesis that it differs from 0.5.

$$\rightarrow \hat{p} = \frac{18}{200} = 0.09.$$

$$H_0: p = p_0 = 0.5$$

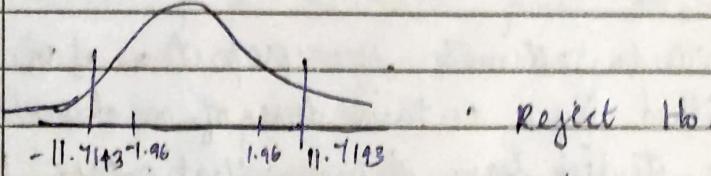
$$H_1: p \neq p_0 \neq 0.5. \quad \text{Two-tailed - "differ".}$$

$$\alpha = 0.05. \quad z = 1.96.$$

$$z = \frac{\hat{p} - p_0}{\text{St}_0(\hat{p})}$$

$$\text{St}_0(\hat{p}) = \sqrt{\frac{p_0(1-p_0)}{n}} = \sqrt{\frac{0.5(1-0.5)}{200}} = \sqrt{\frac{0.25}{200}} = 0.035$$

$$z = \frac{0.09 - 0.035}{0.035} = \frac{0.055}{0.035} = 1.57143.$$



(2x left area $\leftarrow \alpha$) - with p-value from z-table(17).

On Investors in New IT have a historical success rate of 65%. In a recent poll, a random sample of 250 New IT investors showed that only 150 were successful. At $\alpha = 0.01$, is there enough evidence to conclude that the success rate has increased from 65%.

$$\hat{p} = \frac{180}{250} = 0.72 = 72\% \quad \text{Rate} > 65\%$$

$$H_0: p = p_0 = 0.65$$

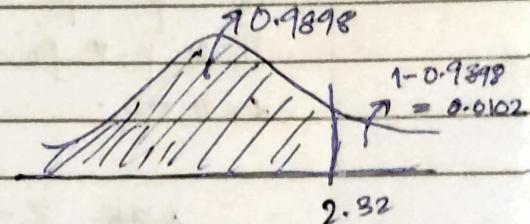
$$H_1: p > p_0 > 0.65 \quad \rightarrow \text{one-tailed}$$

$$\text{SE}(\hat{p}) = \sqrt{\frac{p_0(1-p_0)}{n}} = \sqrt{\frac{0.65(1-0.65)}{250}} \\ = \sqrt{\frac{0.2275}{250}} = 0.0302$$

$$z = \frac{\hat{p} - p_0}{\text{SE}(\hat{p})} = \frac{0.72 - 0.65}{0.0302} = \frac{0.07}{0.0302} = 2.31488$$

$$z = 2.32$$

$$\text{i-value} = 0.9898$$



~~0.9898 > 0.04~~, \therefore fail to reject H_0 .

$$p \geq 1 - 0.9898 = 0.0102$$

$\therefore p \leq 0.04 \quad 0.0102 \leq 0.01$, we reject H_0 .

Qn. Suppose we wish to estimate the proportion of mice that would be killed by a certain dose of a chemotherapy drug. Previous studies have shown that approx 10% of mice are killed by this dose of the drug. Suppose we wish to estimate p to within 0.001 with 99% confidence, how large of a sample size is required?

$$z \cdot \sqrt{\frac{p^*(1-p)}{n}} \leq m$$

$$m = 107 \approx 0.017$$

$$n = 2975$$

$$\approx 2.515 \sqrt{0.01(0.01)} \approx 0.017$$

$$1.2875 \pm 0.017$$

$$n = \sqrt{\frac{1.2875}{0.01}} = 17.35$$

Not sure if correct.

25 Feb

1 Hypothesis Test for comparing two proportions:

- Q1. A study investigated a possible effect of magnetic pulse on the ability of homing pigeons to navigate back to the home loft. Pigeons were randomly selected divided into a magnetic pulse group and a control group. The pigeons were then released from a location 106 km from the home loft.

- 22 out of 38 control group pigeons returned to the home loft. 21 out of the 39 pigeons that received a magnetic pulse returned to the home loft.
- Construct a confidence interval for $p_1 - p_2$
 - Test $H_0: p_1 = p_2$

$$\rightarrow C = \hat{p}_1 - \hat{p}_2 \pm z_{\alpha/2} \times SE(\hat{p}_1 - \hat{p}_2)$$

$$\text{where } SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

$$\hat{p}_1 = \frac{x_1}{n_1} = \frac{22}{38} = 0.58$$

$$z_{\alpha/2} = 1.96$$

$$\hat{p}_2 = \frac{x_2}{n_2} = \frac{21}{39} = 0.54$$

$$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} = 0.113$$

∴

$$CI = 0.09 \pm 1.96(0.113)$$

$$= 0.09 \pm 0.22198$$

To test $H_0: p_1 = p_2$: $= (-0.18198, 0.26148)$

~~$$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$~~

$$H_0: p_1 = p_2 \text{ or } p_1 - p_2 = 0.$$

$$H_1: p_1 \neq p_2 \text{ or } p_1 - p_2 \neq 0.$$

$$z = \frac{\hat{p}_1 - \hat{p}_2}{SE(\hat{p}_1 - \hat{p}_2)}$$

where

$$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$$

$$\text{where } \hat{p} = \frac{x_1 + x_2}{n_1 + n_2}$$

pool sample proportion since our assumption
is $\hat{p}_1 = \hat{p}_2 = \hat{p}$.

$$\hat{p} = \frac{x_1 + x_2}{n_1 + n_2} = \frac{22 + 21}{38 + 39} = 0.558$$

$$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{0.558(1-0.558)\left(\frac{1}{38} + \frac{1}{39}\right)}$$

$$= 0.1132$$

$$\therefore z = \frac{\hat{p}_1 - \hat{p}_2}{SE(\hat{p}_1 - \hat{p}_2)} = \frac{0.58 - 0.59}{0.1132} = -0.353$$

Two-tailed $\alpha = 0.05$.

$$1-\text{Pvalue} = 1 - 0.6368 = 0.3632$$

$$\text{Pvalue} = 0.3632 \times 2 = 0.7264$$

Since $p > \alpha$, we fail to reject H_0 .

~~marks~~

Semester Assignments

Many many assignments.

Indeath the destroyer of assignments.

January ? No No ? February ? Yes Yes ?

- Corrected by Jeshin

1/1

Qn. A machine puts out 16 imperfect articles in a sample of 500. After the machine is overhauled, it puts out 3 imperfect articles in a batch of 100. Has the machine improved?

$$H_0: p_1 = p_2 \quad \text{Unetailed.}$$

$$H_1: p_1 < p_2 \quad ? \quad \text{Left tailed?}$$

$$\hat{p}_1 = \frac{x_1}{n_1} = \frac{16}{500} = 0.032.$$

$$\hat{p}_2 = \frac{x_2}{n_2} = \frac{3}{100} = 0.03.$$

$$Cl = \hat{p}_1 - \hat{p}_2 \pm 2\alpha/2$$

27 Feb.

* Parametric & Non parametric Tests:

Non parametric:

(i) Chi-Square Test: (χ^2).

for goodness of fit.

$$\chi^2 = \sum_{E} (O-E)^2 = \sum_{\text{Expected}} (\text{Observed} - \text{Expected})^2.$$

$$\chi^2 = \sum_{\text{Actual}} (\text{Predicted} - \text{Actual})^2. \quad \chi^2 = \sum_{\text{Predicted}} (\text{Actual} - \text{Predicted})^2$$

χ^2 Assume expected/actual values > 5 .



$\chi^2_{\text{critical}} (\alpha, df_1) \leq \chi^2_{\text{stat}}$. \rightarrow Reject H_0 .

$\chi^2_{\text{critical}} (\alpha, df_1) > \chi^2_{\text{stat}}$. \rightarrow Fail to reject H_0 .

$$df_1 = k-1.$$

Qn. Sample of 100. Allergy in 2015 - Yes: 53%, No: 47%.
2020 - Yes: 49%, No: 51%.

$$H_0: (p_1, p_2) = (53\%, 47\%)$$

$$H_1: (p_1, p_2) \neq (53\%, 47\%).$$

Sample - observed ; population - expected.

$$O_1: 49$$

$$E_1: 53$$

$$O_2: 51$$

$$E_2: 47$$

$$\begin{aligned} \chi^2 &= \frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2} = \frac{(19 - 53)^2}{53} + \frac{(31 - 47)^2}{47} \\ &= 0.302 + 0.39 = \underline{\underline{0.69}}. \end{aligned}$$

$$\chi^2_{\text{critical}} (\alpha, df_1) = \chi^2 (0.05, 2) = \chi^2 (0.05, 1) = 3.891.$$

Since ~~χ^2_{stat}~~ $\chi^2_{\text{critical}} > \chi^2_{\text{stat}}$, we fail to reject H_0 .

p-value: At $df_1 = 1$, $\alpha = 0.9$ $\alpha = 0.1$

$$df_1 = 1: \quad 0.016 \quad \downarrow \quad 0.059706.$$

$$\text{p-value} = \underline{\underline{0.92}}.$$

Since ~~$0.92 > 0.05$~~ $0.92 < 0.05$, we reject H_0 .

Since $0.92 > 0.05$, fail to reject H_0 .

Qn. Some research shows higher number of flight tickets are bought by males in comparison to females with ratio of 2:1. Out of 150 tickets, 88 were bought by males, 62 by females. We need to find out if the experimental manipulation causes the change in the results, or if we are observing a chance variation.

$$2:1 \rightarrow 150.$$

$$\rightarrow 100:50. \quad O_1: 88.$$

$$E_1: E_2. \quad O_2: 62.$$

$$\chi^2 = \frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2} = \frac{(88 - 100)^2}{100} + \frac{(62 - 50)^2}{50} \\ = 1.44 + 2.88 = \underline{\underline{4.32}}.$$

$$df_1 = 2 - 1 = 1.$$

$$\alpha = 0.05.$$

$$\chi^2_{\text{critical}} = (0.05, 1) = 3.841.$$

Since $\chi^2_{\text{critical}} (3.841) \leq \chi^2_{\text{stat}} (4.32)$, we reject H_0 .

* Chi-Square test of Independence:

Qn. $n = 8100$.

	Cancer	No cancer	Total
Smokers	60.0%	300.0%	360
No smokers	40.0%	1200.0%	1240
% smokers	60%	10%	
Total	100	8000	8100

$$E_{11} = \frac{\text{Row}_1}{N} \times \frac{\text{Col}_1}{N} \times N = \frac{860}{8100} \times \frac{100}{8100} \times 8100 = \underline{10.62}$$

$$E_{12} = \frac{\text{Row}_1}{N} \times \frac{\text{Col}_2}{N} \times N = \frac{860}{8100} \times \frac{8000}{8100} \times 8100 = \underline{849.4}$$

$$E_{21} = \frac{\text{Row}_2}{N} \times \frac{\text{Col}_1}{N} \times N = \frac{7240}{8100} \times \frac{100}{8100} \times 8100 = \underline{89.4}$$

$$E_{22} = \frac{\text{Row}_2}{N} \times \frac{\text{Col}_2}{N} \times N = \frac{7240}{8100} \times \frac{8000}{8100} \times 8100 = \underline{7150.62}$$

$$\begin{aligned} \chi^2 &= \frac{(O_{11} - E_{11})^2}{E_{11}} + \frac{(O_{12} - E_{12})^2}{E_{12}} + \frac{(O_{21} - E_{21})^2}{E_{21}} + \frac{(O_{22} - E_{22})^2}{E_{22}} \\ &= \frac{(60 - 10.62)^2}{10.62} + \frac{(800 - 849.4)^2}{849.4} + \frac{(40 - 89.4)^2}{89.4} + \frac{(7200 - 7150.62)^2}{7150.62} \\ &= 229.6 + 2873 + 27.3 + 0.341 \\ &= \underline{260.119}. \end{aligned}$$

$$df = (\text{Row}-1)(\text{Col}-1) = (2-1)(2-1) = 1 \cdot 1 = 1.$$

$$\chi^2_{\text{critical}}(0.05, 1) = \underline{3.841}.$$

Since $\chi^2_{\text{critical}}(0.05, 1) \leq \chi^2_{\text{stat}}(260.119)$, we reject H_0 .

28 Feb.

11

* K-S test (Kolmogorov-Smirnov):

To compare sample & population distribution.
Observed & theoretical.

Qn. In a study done from various streams of a college of 60 students, with equal no. of students drawn from each stream, the intention of the students to join the Adventure club of college noted after the interviewed are listed below.

	BSc	BA	BCom	MA	MCom
# in class	5	9	11	16	19

It was expected that 12 students from each class would join the Adventure club. Using K-S, find if there is any difference among student classes with regards to their intention of joining the club.

$$\rightarrow H_0: F_{\text{O}(x)} = F_{\text{N}(x)} \rightarrow \text{no difference.}$$
$$H_1: F_{\text{O}(x)} \neq F_{\text{N}(x)} \rightarrow \text{Yes difference.}$$

	Actual/Obs	Theoretical	$F_{\text{O}(x)}$ CDF	$F_{\text{N}(x)}$ CDF	$D = F_{\text{O}(x)} - F_{\text{N}(x)} $
BSc	5	12	5/60	12/60	7/60
BA	9	12	14/60	24/60	10/60
BCom	11	12	25/60	36/60	11/60
MA	16	12	41/60	48/60	7/60
MCom	19	12	60/60	60/60	0
Total	60	60			

$$D_{\text{stat}} = \max |F_{\text{O}(x)} - F_{\text{N}(x)}| = 11/60 = 0.183.$$

$$n > 10 \rightarrow \alpha = 0.05 \rightarrow D_{\text{critical}} = \frac{1.36}{\sqrt{60}} = 0.175.$$

Since $D_{\text{stat}} (0.183) > D_{\text{critical}} (0.175)$,

reject H_0 .

Qn. Consider the data points

1.91 0.26 1.91 0.33 0.55 0.77 1.46 +94.1.12

Is there any evidence to suggest that the data were not randomly sample from a uniform (0,2) distribution?

$$\rightarrow H_0: F(x) = F_0(x).$$

$$H_1: F(x) \neq F_0(x).$$

* K-S - Two sample test:
Comparing distribution of 2 samples.

$$D = \max |F_{n_1}(x) - F_{n_2}(x)|.$$

Small sample : when $n_1 + n_2 < 10$.

Large sample : when n_1 and/or $n_2 \geq 10$.

Ex. consider a survey on two different universities of the PG students on the topic of their willingness to join the Research Funding Project on AI. The following ~~are~~ is the results obtained:

Uni 1	3	2	3	5	8	9	8	8
Uni 2	2	8	2	4	4	3	6	

Determine whether the samples for uni 1 & 2 come from the same distribution.

- $\rightarrow H_0: F_1(x) = F_2(x)$ \rightarrow Same distribution.
 $H_1: F_1(x) \neq F_2(x).$ \rightarrow Different distribution.

Frequency

Items	Uni 1	Uni 2	CDF $F_1(x)$	CDF $F_2(x)$	$D = F_1 - F_2 $
2	1	2	1/8	2/4	• 0.161.
3	2	1	3/8	3/4	• 0.0536.
4	0	2	3/8	5/4	0.34.
5	1	0	4/8	5/4	0.214.
6	0	1	4/8	6/4	→ 0.36.
8	3	1	7/8	7/4	0.125.
9	1	0	8/8	7/4	0

Total: 8 7

$\therefore D = 0.36;$ $D_{\text{critical}}(0.05, 8, 7) = 40/56 = 0.7143.$ $D_{\text{test}} < D_{\text{critical}},$ ~~so~~ fail to reject.

* Sign Test:

- Non parametric test.
- Using population median \bar{N} .
- If $T_{\text{stat}} > T_c$, fail to reject H_0 ; $T_{\text{stat}} < T_{c*}$, reject H_0 .

$$H_0: \bar{N} = N_0$$

H_F

$$H_1: \bar{N} \neq N_0 \quad \leftarrow \text{Two tailed.}$$

$$H_1: \bar{N} < N_0 \quad \leftarrow \text{left tailed.}$$

$$H_1: \bar{N} > N_0 \quad \leftarrow \text{right tailed.}$$

Qn. Following are the responses to the question "How many hours does than study before a major statistic test, mind laid?"

6 5 1 2 2 5 7 5 8 7 4 7.

Use the sign test to test the hypothesis at 5% sig. lvl. that the median # hours a student studies before a test is 3. Given that the critical value of sign test for $n=11$ at 0.05α for two tailed is 1.

→ $n \leq 25$, ∴ small sample.

$$H_0: \bar{N} = 3.$$

$$T_c = 1. \quad n = 11.$$

$$H_1: \bar{N} \neq 3 \quad \text{Two-tailed.}$$

$$N_0 = 3.$$

6 - 3	5 - 3	1 - 3	2 - 3	2 - 3	5 - 3	7 - 3	5 - 3	3 - 3	7 - 3	4 - 3	7 - 3
3	2	- 2	- 1	- 1	2	4	2	0	4	1	4
+	+	-	-	-	+	+	+	0	+	+	+

$$T^+ = 8. \quad ; \quad T^- = 3. \quad ; \quad \text{Rej. H}_0.$$

$$\Rightarrow T = \min(T^+, T^-) = \min(8, 3) = 3.$$

$\xrightarrow{\text{Reject } H_0}$ T_c

Since $T(3) > T_c(1)$, we fail to reject.

6 March 5th Module:

* Statistical Decision Theory:

Complete Class of Decision rules:

Set of decision rules with less error.

Minimal Complete Class of Decision rules:

'just enough' sort of set, no removal of rules needed, no addition either. just enough.

Loss function; Risk function.

Minimal EEDR: Smallest subset of the complete class that contains all the decision rules that are optimal with respect(?) to the loss function.