

15<sup>th</sup> February 2023

Practical Day 1

$\chi^2$  distribution

Let  $X_1, X_2, \dots, X_n$  be ~~n~~ mutually independent standard normal ~~variables~~ random variables. Then the sum of squares

$$U = \sum_{i=1}^n (X_i^2)$$

is said to follow  $\chi^2$  distribution

with  $n$  degrees of freedom.

It has the following pdf

$$f(u) = \begin{cases} \frac{1}{2^{n/2}} \frac{1}{\Gamma(n/2)} [e^{-u/2} u^{n/2 - 1}] & 0 < u < \infty \\ 0 & \text{otherwise} \end{cases}$$

The distribution is characterized by its degrees of freedom ( $n$ ).

We will write it as  $U \sim \chi_n^2$

Let  $X$  be a standard normal random variable &  $Y$  be another random variable following "chi-sq dist" with degrees of freedom  $n$ . Furthermore,  $X$  &  $Y$  are independent. Then the new variable  $t = \frac{X}{\sqrt{\frac{Y}{n}}}$  is

said to follow t-distribution, with  $n$  degrees of freedom.

$t$  has the following pdf

$$f(t) = \begin{cases} \frac{1}{\sqrt{n} B(\frac{n}{2}, \frac{1}{2})} \left(1 + \frac{t^2}{n}\right)^{-\frac{n+1}{2}} & t \in \mathbb{R} \\ 0 & \text{otherwise} \end{cases}$$

We write  $t$  as  $t \sim tn$

Q2) For  $n = 10, 20, 30, \dots, 60$  plot the pdf of a  $t_n$  distribution  
in separate panels of the same graph & comment.

$c_1 := t$  [-5 to 5 with step value 0.01  
make patterned data]

$c_2 := f(t)$  ( $n=10$ )

calc → Prob Dist →  $t$

$c_3 := f(t)$  ( $n=20$ )

Prob density

( $n=60$ )

Input →  $t$

~~Output~~ optional storage →  $c_2$

We notice that as  $n \rightarrow \infty$ ,  $t$ -dist tends to be normal dist.

Graph: same as chisq,

[Apply overlaid graphs to understand skewness and kurtosis]

Skewness → Symmetric

Area under tail decreasing with increasing  $n$ .

► Mesokurtic

Kurtosis Area under tail / heavy or thick tail  
As  $k \downarrow$ ,  $\rightarrow$  mesokurtic

Area under tail ↑ Heavy / thick tail ↑

Q3) Plot the pdf of the  $t$ -distribution having 10 df and the "std. ~~density~~ normal dist" on the same graph and comment.

[-4 to 4 gap 0.01]

Kurtosis more for  $t$ -dist.

→ F-distribution

let  $X$  and  $Y$  be two independently distributed  $\chi^2$  random variables with degrees of freedom  $m$  &  $n$  respectively. Then the new random variable

$F = \frac{(x/m)}{(v/n)}$  is said to follow F-distribution with  $(m, n)$  degrees of freedom. It has the following pdf.

$$F(f) = \begin{cases} \frac{(m/n)^{m/2}}{B(m/2, n/2)} f^{(m-2)/2} (1 + \frac{m}{n}f)^{-(m+n)/2} & 0 < f < \infty \\ 0 & \text{otherwise} \end{cases}$$

We write it as  $F \sim F_{m,n}$

Q1) Plot the pdf of the F distribution having df  $(n_1, n_2)$  for the following choices of  $n_1$  &  $n_2$

i)  $n_1 = 10$  &  $n_2 = 10, 20, 30, 40$

ii)  $n_1 = 20$  &  $n_2 = 10, 20, 30, 40$

iii)  $n_1 = 30$  &  $n_2 = 10, 20, 30, 40$

iv)  $n_1 = 40$  &  $n_2 = 10, 20, 30, 40$

+ve skewed

Kurtosis decrease

$c_1 = f$

$c_2 = n_2 = 10$

1st  
March  
2023

Tests & Confidence intervals related  
to univariate normal pop<sup>n</sup>. (AI)

Practical Day 3

Q5) Nine items of a sample had the following values 45, 47, 50, 52, 48, 47, 49, 53 & 51. Does the mean of the normal population from which these 9 items are drawn differ significantly from an assumed pop<sup>n</sup> mean of 47.5 (i) when pop<sup>n</sup> std is unknown & given to be 2.5 and (ii) when pop<sup>n</sup> std is unknown. For case (ii), obtain 95.7% confidence interval for the pop<sup>n</sup> mean.

Let  $X$  be a random variable denoting the value of a randomly selected item.

Given that  $X \sim N(\mu, \sigma^2)$ ; Sample size:  $n = 9$ .

Let  $x_1, x_2, \dots, x_n$  be a random sample from the dist<sup>n</sup> of  $X$ . We want to test  $H_0: \mu = 47.5$

against  $H_1: \mu \neq 47.5$

~~Since  $\sigma$  is unknown, the apt test is one sample t-test.~~

ii) Since  $\sigma$  is unknown. Here, the apt test is one sample t-test. Here the test statistic is given by  $t = \frac{\sqrt{n}(\bar{x} - \mu_0)}{s}$  where  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

then  $t \sim t_{n-1}$ .

$|t_{obs}| > t_{\alpha/2, n-1}$ .  $t$  denotes upper  $(\alpha/2)^{th}$  quantile of  $t$ -dist<sup>n</sup> with  $(n-1)$  df.

Calc  $\rightarrow$  std deviation (Input c1)

$$s = 2.61937 \text{ and } t = \frac{\sqrt{n}(\bar{x} - \mu_0)}{s} = 1.84523$$

$$\alpha = 0.05 \text{ and } \alpha/2 = 0.025$$

$$P(t > t_{\alpha/2, n-1}) = \alpha/2 \Rightarrow P(t < t_{\alpha/2, n-1}) = 0.975$$

$$\Rightarrow t_{\alpha/2, n-1} = P_t^{-1}(0.975) \\ = 2.306$$

$$\text{and } |t_{obs}| = 1.84523$$

$\therefore |t_{obs}| < t_{\alpha/2, n-1}$  is what we get.

$\therefore$  On the basis of the given sample we do not reject  $H_0$  at 5% level of significance.

$\therefore$  We can conclude that the pop' mean is not significantly different from 47.5

i) Test of  $\mu$  when  $\sigma^2$  is known.

ii) Given that  $\sigma = 2.5$ . Since  $\sigma$  is known, the apt test should be 1-sample Z-test.

Here the test statistic is given by

$$Z = \frac{\sqrt{n}(\bar{x} - \mu_0)}{\sigma} \sim N(0, 1)$$

Both-tailed test

We will reject  $H_0$  at level  $\alpha$  if  $|Z_{obs}| > Z_{\alpha/2}$ , where  $Z_{\alpha/2}$  denotes the upper  $(\alpha/2)^{th}$  quantile of  $N(0, 1)$

~~sample~~ → 1 set of sample  
guarantee participant no.

2-sample → 2 set of sample

Calculation:  $n = 9$ ;  $\mu_0 = 47.5$ ;  $\sigma = 2.5$ ;  $\bar{X} = \frac{1}{n} \sum X_i$

$$Z = \frac{\sqrt{n}(\bar{X} - \mu_0)}{\sigma} = 1.932 \quad \leftarrow$$

using Minitab

$$\alpha = 1 - 0.95 = 0.05$$

$$\therefore Z_{\alpha/2} = Z_{0.025}$$

[Calc → Prob dist<sup>n</sup> → Normal → Inverse cumulative  
 $P(Z > Z_{\alpha/2}) = \alpha/2 \rightarrow P(Z < Z_{\alpha/2}) = 1 - \alpha/2 = 0.975$

$$\Rightarrow Z_{\alpha/2} = \Phi_z^{-1}(0.975)$$

$$= 1.959$$

∴ [Input constant → 0.975 Optional storage : blank]

$$\therefore Z_{\text{obs}} = 1.933 \quad \therefore (Z_{\text{obs}}) = 1.933 < 1.959 = Z_{0.025}$$

Hence, on the basis of the given sample, we do not reject  $H_0$  at 5% level of significance

∴ We can conclude that the population mean is not significantly different from 47.8.

15<sup>th</sup> March 2023 (A1) (A2)

Practical Day 4

- g6) A firm manufacturing rivets to limit variation in their lengths as much as possible. The length (in cms) of 10 rivets manufactured by a new process is 2.10, 1.94, 2.0, 2.07, 2.12, 1.96, 1.93, 1.98, 2.20, 1.88. In the past the std deviation of length of rivets by the firm has been 0.145 cms. Examine whether the new process seems to be ~~more~~ cm ~~more~~ superior to the old in case
- the mean length is known to be 2.00 cm and
  - the mean length is unknown. Obtain 95% confidence interval for the std deviation when the mean length is unknown. Assume the length of rivets to be normally distributed.

Sol":

let  $X$  be a random variable denoting the length<sup>(in cm)</sup> of a randomly selected rivet manufactured by the new process of the firm.

In the past, the standard deviation of the rivets manufactured by the firm has been  $0.195 \text{ cm} = \sigma$  (say)

let  $X_1, X_2, \dots, X_{10} \sim X$  be a random sample of size 10.

We are given a random sample of size 10. The sample values are 2.10, 1.91, 2.00, 2.07, 2.12, 1.96, 1.93, 1.98, 2.20, 1.88.

Given mean length is 2.00 cm. That is the population mean  $\mu$  is given.

We will be testing  $\sigma^2$  by the following hypothesis  $H_0: \sigma^2 = 0.145$  against  $H_1: \sigma^2 < 0.145$

Reason behind the formulation of such an hypothesis : We need to check if the new process (whose sample values are given) is superior than the old or pre-existing process of manufacturing by the firm. which means superiority implies less variation and hence the hypothesis

(a) Under  $H_0$ ,

Consider the test statistic  $S^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_0)^2$

$$\text{let } T = \sum_{i=1}^n (x_i - \mu_0)^2 = (ns^2 / \sigma_0^2)$$

Clearly,  $T \sim \chi_n^2$ . of the rejection rule is as follows, know

We will reject  $H_0$  if  $s^2 < 0.195^2$ , ie a lower value of  $\frac{s^2}{\sigma_0^2}$  less than 1 will indicate departure from  $H_0$ .

Given level of significance is  $\alpha = 0.05$ .

Thus  $P(\text{Type I error}) = \alpha \Rightarrow P(\text{Reject } H_0 | H_0 \text{ is true}) = \alpha$

$$\Rightarrow P(T < k_1 | \sigma_0 = 0.195) = \alpha$$

$$\Rightarrow k_1 = F_{\chi_n^2}^{-1}(\alpha) = \chi_{n; 1-\alpha}^2 \text{ the lower } \alpha\text{-point}$$

Hence the test rule is reject  $H_0$  at a level of significance iff  $T_{\text{obs}} < \chi_{n; 1-\alpha}^2$

$$\text{We see } T_{\text{obs}} = \frac{n s^2}{\sigma_0^2} = 8^2 = 0.00942$$

$$T_{\text{obs}} = 4.48038$$

and using minitab we find that  $\chi_{n; 1-\alpha}^2 = 3.94030$

Hence we see  $T_{\text{obs}} = 4.48 > 3.9403 = \chi_{n; 1-\alpha}^2$   
So we do not reject  $H_0$  at a level of significance

$\therefore H_0$  is accepted when  $\mu$  is 2 cm.

Hence, the new process of manufacture do not differ significantly with respect to the old one in terms of superiority.

(b) Here the mean length is unknown.  
 Consider  $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$  where  $\bar{x}$  = sample mean.

Now consider the test statistic  $T = \frac{(n-1)s^2}{\sigma_0^2}$  under  $H_0$

Then  $T = \frac{\sum (x_i - \bar{x})^2}{\sigma_0^2} \Rightarrow T \sim \chi_{n-1}^2$  where  $n = 10$   
 Under  $H_0$ ,  $T \sim \chi_{n-1}^2$

Now  $P_{H_0} [T_0 < c] = \alpha \Rightarrow$  ~~the~~  $c$  is the ~~upper~~ lower alpha point

Hence we will reject  $H_0$  iff at  $\alpha$  level of significance  
 iff  $T_{obs} < \chi_{n-1, 1-\alpha}^2$

We see  $T_{obs} = \frac{(n-1)s^2}{\sigma_0^2}$ ; Using minitab, we find

Now  $\chi_{n-1, 1-\alpha}^2 = 3.32511$ .  
 $(n-1)s^2 = 0.0909$  and so  $T_{obs} = 4.3262$ .

we find that  $T_{obs} = 4.32 \neq 3.32 = \chi_{n-1, \alpha}^2$ .  
 Thus, we do not reject  $H_0$  at a level of significance

$\therefore H_0$  is accepted when mean length is unknown.  
 Thus, the new process of manufacture do not differ significantly with respect to the old one in terms of superiority.

let  $D$  be a dist. Hence  $D_{\alpha}$  means it is a point such that  $P_D(D_{\alpha}) = \alpha$ .

ज्ञानी यज्ञम् इवं कुलाची-ग्रन्थ

अत इति '६' वर्षस्तु रूपम् विद्या यज्ञम् इति अनुवादः

representation द्वारा  $D_{\alpha} = \text{quantile}(D, 1-\alpha)$

Example.  $K = \chi_{n-1, 1-\alpha}^2$  means  $P_{\chi_{n-1}^2}(\chi_{n-1, 1-\alpha}^2 > K) = 1-\alpha$

$\Rightarrow 1 - P_{\chi_{n-1}^2}(\chi_{n-1, 1-\alpha}^2 < K) = 1-\alpha \Rightarrow P_{\chi_{n-1}^2}(\chi_{n-1, 1-\alpha}^2 < K) = P_{\chi_{n-1}^2}^{-1}(1-\alpha)$

$\Rightarrow \chi_{n-1, 1-\alpha}^2 = P_{\chi_{n-1}^2}^{-1}(1-\alpha)$

Under null hypothesis  $H_0$ , consider the test statistic  $T$  defined as

$$\frac{\bar{X} - \bar{Y}}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} = \frac{\bar{X} - \bar{Y}}{S \sqrt{\frac{1}{13} + \frac{1}{12}}}$$

where

$$S^2 = \frac{(n_1-1) s_1^2 + (n_2-1) s_2^2}{(n_1+n_2-2)} \quad \text{where } s_1^2 = \frac{1}{n_1-1} \sum_{i=1}^{n_1} (x_i - \bar{x})^2$$

and  $s_2^2 = \frac{1}{n_2-1} \sum_{j=1}^{n_2} (\bar{y}_j - \bar{Y})^2$ . where  $\bar{x}, \bar{Y}$  are sample means.

Now we know  $\frac{(n_1-1) s_1^2}{\sigma^2} \sim \chi_{n_1-1}^2$  and  $\frac{(n_2-1) s_2^2}{\sigma^2} \sim \chi_{n_2-1}^2$

$$\text{So } (n_1+n_2-2) S^2 / \sigma^2 \sim \chi_{n_1-1}^2 + \chi_{n_2-1}^2 = \chi_{n_1+n_2-2}^2$$

Now  $\bar{X} \sim N(\mu_1, \frac{\sigma^2}{n_1})$  and  $\bar{Y} \sim N(\mu_2, \frac{\sigma^2}{n_2})$  under  $H_0$

$$\Rightarrow (\bar{X} - \bar{Y}) \sim N(\mu_1 - \mu_2, \sigma^2(\frac{1}{n_1} + \frac{1}{n_2})) \text{ under } H_0$$

$$= N(0, \sigma^2(\frac{1}{n_1} + \frac{1}{n_2}))$$

$$\Rightarrow \frac{\bar{X} - \bar{Y}}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim N(0, 1).$$

$$\text{and } (n_1+n_2-2) S^2 / \sigma^2 \sim \chi_{n_1+n_2-2}^2$$

$$\therefore T = \frac{\bar{X} - \bar{Y}}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \Rightarrow \frac{(\bar{X} - \bar{Y})}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim \sqrt{\frac{S^2(n_1+n_2-2)}{\sigma^2} \left( \frac{1}{\frac{1}{n_1} + \frac{1}{n_2}} \right)}$$

$$\Rightarrow T = \frac{N(0, 1)}{\sqrt{\frac{\chi_{n_1+n_2-2}^2}{n_1+n_2-2}}} \sim t_{n_1+n_2-2}$$

$$\Rightarrow T \sim t_{n_1+n_2-2} \text{ under } H_0$$

so accordingly, we will be rejecting  $H_0$  iff we get  $T_{obs} < t_{-\alpha; n_1+n_2-2}$  at  $\alpha = 0.05$  level of significance.

Using 2-sample t-test in minitab, we find

$$T_{obs} = -0.89 \quad \text{and value of } S = 8.3805$$

$$\text{and } t_{\alpha/2; n_1+n_2-2} = -1.71 \quad [\text{Here } \alpha = 0.05, n_1 = 13, n_2 = 12]$$

Thus, we see  $T_{obs} = -0.89 \not< -1.71 = t_{\alpha/2; n_1+n_2-2}$

Conclusion: Thus we do not reject  $H_0$  at  $\alpha$  level of significance. Hence, the average daily milk production (in pounds) do not differ significantly between cows fed with field watered alfalfa and dewarfed alfalfa.

## CI types

CI for  $\mu$  ( $\sigma$  unknown)

: 1-sample T (using minitab)

CI for  $\mu$  ( $\sigma$  known)

: 1-sample Z (using minitab)

CI for  $\sigma$

( $\mu$  known) :  $\chi_n^2$  (manually)

CI for  $\sigma$  ( $\mu$  unknown) :  $\chi_{n-1}^2$  (manually)

CI for  $\mu_1 - \mu_2$  ( $\sigma_1 = \sigma_2$  <sup>unknown</sup>) : 2-sample T (using minitab)

CI for  $\sigma_1/\sigma_2$  ( $\mu_1, \mu_2$  unknown) : 2-variances (using minitab)  
(Sample std ~~will be provided~~ will be provided)

CI for  $\sigma_1/\sigma_2$  ( $\mu_1, \mu_2$  known) :  $F_{n_1, n_2}$  (manually)

$$\left( \frac{s_1^2}{s_2^2} \left( F_{\alpha/2; n_1, n_2} \right)^{-1}, \frac{s_1^2}{s_2^2} \left( F_{1-\alpha/2; n_1, n_2} \right)^{-1} \right).$$

B2b

Q) It is known ~~o~~ that the mean diameter of rivets produced by two factories is practically the same but std deviation may differ. For 22 rivets produced by factory I, the std deviation is 2.9 mm, while for 16 rivets produced by factory II the s.d. is 3.8 mm. Do you think that the products of factory I are of better quality than that of factory II? Also obtain 95% confidence interval

Sdn:

let  $X_1$  denote the diameter of a randomly selected rivet produced by factory I.  
 $Y_2$  denote the diameter of a randomly selected rivet produced by factory II.  
 $X_1, Y_2$  are independent as the factories are different.

Let  $X_1 \sim N(\mu, \sigma_1^2)$  and  $Y_2 \sim N(\mu, \sigma_2^2)$  where  $\mu, \sigma_1, \sigma_2$  are unknown.

We are to test for ~~the~~  $\sigma_1^2 / \sigma_2^2$  the population variance ratio. For better quality, the factory having lesser variance will be preferred.

So let us test  $H_0: \frac{\sigma_1^2}{\sigma_2^2} = 1$  vs  $H_1: \frac{\sigma_1^2}{\sigma_2^2} \neq 1 (\neq 1)$

Now let  $X_1, \dots, X_{22}$  be a random sample of size 22 taken from the distribution of  $X$

and let  $Y_1, \dots, Y_{16}$  be a random sample of size 16 taken from the distribution of  $Y$ .

From given problem we know sample standard deviation of following  $Y$ -distribution is

$$s_1 = 2.9$$

and that following  $Y$ -distribution is  $s_2 = 3.8$

consider the test statistic  $F_1 = \frac{s_{x_0}^2}{s_{y_0}^2}$  where

$$s_{x_0}^2 = \frac{1}{n_1-1} \sum_{i=1}^{n_1} (x_i - \bar{x})^2 \quad \text{and}$$

$$s_{y_0}^2 = \frac{1}{n_2-1} \sum_{i=1}^{n_2} (y_i - \bar{y})^2$$

$$\text{where } s_0 \quad F = \frac{\frac{(n_1-1)s_{x_0}^2}{(n_1-1)s_{y_0}^2}}{\frac{(n_2-1)s_{y_0}^2}{(n_2-1)s_{x_0}^2}} = \frac{s_{x_0}^2}{s_{y_0}^2} \cdot \frac{\frac{s_1^2}{s_2^2}}{\frac{s_2^2}{s_1^2}} = \frac{s_{x_0}^2}{s_{y_0}^2} \cdot \frac{s_1^2}{s_2^2}$$

Hence

$$F \sim F_{n_1-1, n_2-1}$$

Small value of  $T$  goes in favour of  $H_1$ . Hence, a left-tailed test based on  $T$  will be appropriate.

Rejection Rule [We reject  $H_0$  if  $F_{\text{obs}} < F_{1-\alpha; n_1-1, n_2-1}$  where

$F_{n_1-1, n_2-1; 1-\alpha}$  denotes lower  $\alpha^{\text{th}}$  point of  $F_{n_1-1, n_2-1}$  distribution.

Using minitab, we found  $F_{\text{obs}} = 0.58$ .

and  $F_{1-\alpha; n_1-1, n_2-1}$

Hence we fail to reject  $H_0$  since we do not have significant evidence to conclude that products of factory I is of better quality than factory II.

Methods of estimation

- MME** 1. In a sample of size 8 from the exponential pop" with specification  $f(x) = \frac{1}{\theta} \exp(-\frac{1}{\theta}(x-\theta)) : \theta \leq x < \infty$ , the observations were 3.71, 9.36, 16.21, 2.67, 6.78, 8.92, 10.22, 4.31. Give the moment estimate of  $\theta$  from sample.

In this problem, we are given a random sample of size 8 drawn from an exponential pop" with specification

$$f(x) = \frac{1}{\theta} \exp(-\frac{1}{\theta}(x-\theta)) \quad \theta \leq x < \infty$$

We are to obtain the moment estimate of  $\theta$  from the sample  $x_1, x_2, \dots, x_n$  drawn

In method of moments, we get from dist' given by  $f(x; \theta)$  equate the sample moments with the corresponding pop" moments. If there are  $k$  unknown parameters to be estimated then generally the 1st raw moment about zero and the second to the  $k^{\text{th}}$  central moments are equated.

Since here, we have to estimate  $\theta$ , hence by equating the sample mean with the pop" mean we get the eqn.

(1st order) pop" mean  $\leftarrow \mu'_1 = \bar{x}$  → sample mean where  $\mu'_1 = E(X)$ ,  $X$  following exp dist' characterized by given pdf  $f(x)$  and

$$\bar{x} = \frac{1}{n} \left( \sum_{i=1}^n x_i \right) \text{ for a given sample.}$$

$x_1, \dots, x_n$  drawn from the dist' of  $X$ .

$$\begin{aligned}
 \text{Now } E(X) &= \frac{1}{2} \int_0^{\infty} x e^{-\frac{1}{2}(x-\theta)} dx = \frac{1}{2} x \int_0^{\infty} e^{-\frac{1}{2}(x-\theta)} dx \\
 &= \left[ \frac{1}{2} x \cdot \frac{e^{-\frac{1}{2}(x-\theta)}}{(-\frac{1}{2})} \right]_0^{\infty} + \left[ \frac{e^{-\frac{1}{2}(x-\theta)}}{(-\frac{1}{2})} \right]_0^{\infty} \\
 &= \boxed{0} + \left( \frac{1}{2} \right) = \left( \theta + \left( \frac{1}{2} \right) \right)
 \end{aligned}$$

$$\text{Now } E(X) = \bar{x} \text{ so } \theta + 2 = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\theta}_{MME} = \frac{1}{n} \left( \sum_{i=1}^n x_i \right) - 2$$

Now from the given sample data, the sample mean is  $\bar{x} = 7.77625$

Hence  $\hat{\theta}_{MME} = \bar{x} - 2 = 5.77$  which is the required moment estimate of the  $\theta$  from the sample.

**Q2)** In drawing 25 balls with replacement from an urn containing white and black balls. If 11 are found to be white. Find the maximum likelihood estimate of  $P$ , the proportion of black balls in the urn. It is known that  $P$  can assume only one of the following values 0.3, 0.35, 0.40, 0.45, 0.50, 0.55, 0.60, 0.65, 0.70. What would be MLE of  $P$  in this case?

Let  $X$  be a random variable denoting the no of black balls obtained out of 25 balls drawn randomly from an urn with replacement, containing white and black balls.

So  $X \sim \text{Bin}(25, p)$ , where  $p$  is the proportion of black balls in an urn.

Here it is given that out of 25 balls drawn, 11 are found to be white, so 14 are black.

Now we have only one sample where there are 14 black balls.

$$\text{Hence, } P(X=14) = \binom{25}{14} p^{14} (1-p)^{11}$$

$$\text{Now } L(p) = \prod_{i=1}^{14} f(p; x_i) = P(X=14)$$

is the likelihood function which is to be checked for all different values of  $p$  given in the question.

[Construct a table :  $p$  vs  $L(p)$ ]

$$\begin{array}{|c|c|} \hline p & L(p) \\ \hline \end{array}$$

Using minitab, the table of values are noted.

We find that  $L(p)$  takes maximum value when  $p = 0.55$ .

So the required maximum likelihood estimate

$$\text{is } \hat{p} = 0.55$$

Q3) The length of life recorded in hours for 10 electron tubes, were 980, 1020, 995, 1015, 990, 1030, 975, 950, 1050, 870. Assume the lifetime is distributed in the form :  $f(x) = \frac{1}{\theta} \exp\left(-\frac{x}{\theta}\right); \theta > 0$

$0 < x < \infty$ . Obtain the MLE of  $\theta$ . Also estimate the probability that the electron tube will survive at least 100 hours.

MLE

let  $X$  be a random variable denoting the length of life (in hours) of a randomly selected electron tube.

Here  $X \sim \text{Exp}(\text{mean} = \theta)$  given by the

$$\text{pdf } f(x) = \frac{1}{\theta} e^{-\frac{x}{\theta}} \quad \theta > 0, \quad 0 < x < \infty$$

Here we are given the length of life for ~~n~~ ten electron tubes. We have to obtain the MLE of  $\theta$ .

We know that for an observed sample

$(x) = (x_1, x_2, \dots, x_n)$  drawn from the given dist<sup>n</sup>, the likelihood function is given by  ~~$L(\theta | x)$~~

$$\text{given by } L(\theta | x) = \prod_{i=1}^n f(\theta; x_i)$$

$$= \frac{1}{\theta^n} e^{-\frac{1}{\theta} \sum x_i}$$

Hence the log-likelihood function is given by

$$L = \log L = -n \log \theta - \frac{1}{\theta} \sum x_i = \varphi$$

likelihood Diff w.r.t.  $\theta$

$$\text{equation. } \Rightarrow -\frac{n}{\theta} + \frac{1}{\theta^2} \sum x_i = 0$$

$$\Rightarrow n = \frac{1}{\theta} \sum_{i=1}^n (x_i)$$

$$\Rightarrow \theta = \left( \frac{\sum x_i}{n} \right). \text{ Now check the second derivative criteria}$$

$$L'(\theta) = \left( -\frac{n}{\theta^2} - \frac{2}{\theta^3} \sum x_i \right) < 0 \text{ so we get}$$

$\theta = \left( \frac{1}{n} \sum x_i \right)$  is the value for which  $L(\theta)$  is maximum

Thus,  $\hat{\theta}_{MLE} = \frac{1}{n} \left( \sum x_i \right)$  as the maximum likelihood estimator of  $\theta$

capital e that the  $x_i$ 's are

and  $\hat{\theta}$  is a test statistic in general.

Hence the ML estimate is (calculate)

$$\hat{\theta}_{MLE} = 987.5$$

91) The following table gives the frequency dist<sup>n</sup> of the number of red blood corpuscles (RBC) per cell of a hemocytometer. Obtain the MLE for the true average no of RBC per cell using these data. Also obtain the MLE of the probability that no of RBC per cell is at most one

RBC count	0	1	2	3	4	5	+
Cells count	143	156	68	27	5	1	100

Let  $X$  be a random variable denoting the no of red blood corpuscles in a randomly chosen cell of a hemocytometer.

$\therefore X \sim \text{Poisson}(\lambda)$  where  $\lambda$  is unknown, follows a pdf  $f(x) = \frac{e^{-\lambda} \lambda^x}{x!}$ ,  $x=0, 1, 2, \dots$   
We are to obtain MLE of  $\lambda$ . and  $\lambda > 0$ .

So far an observed sample as given, there are total 100 observations, and the sample data is  $\bar{x} = (x_1, x_2, \dots, x_{100})$ . So the likelihood function is given by  $L(\lambda | \bar{x}) = \prod_{i=1}^{100} f(\lambda | x_i) = (e^{-100\lambda} \lambda^{\sum x_i}) / (\prod_{i=1}^{100} x_i!)$

Hence the log likelihood function is given by  $L_1 = \log L$

$$= -100\lambda + \sum_{i=1}^{100} x_i \log \lambda - \log(A) \quad \text{where } A \text{ is the denominator}$$

$\therefore$  Differentiating  $L_1$  wrt  $\lambda$  we get the critical pt when of  $L(\lambda | \bar{x})$  equated with zero

$$-100 + \frac{1}{\lambda} \sum_{i=1}^{100} x_i \Rightarrow \lambda = \left( \frac{100}{\sum_{i=1}^{100} x_i} \right)^{-1}$$

Now double differentiating  $L_1$  wrt  $\lambda$  we get  $L_1'' = -\frac{1}{2} \sum x_i$   
 $\Rightarrow L_1''|_{\lambda = \frac{1}{\sum x_i}} < 0$  so clearly at  $\lambda = \frac{1}{\sum x_i}$ ,

$L_1$  is maximum and so is  $L$ . Thus we get the maximum likelihood estimator  $\hat{\lambda} = \frac{1}{100} \sum_{i=1}^{100} x_i$ . So, now using the available sample data, we get the maximum likelihood estimate  $\lambda = \frac{1}{100} \sum_{i=1}^{100} x_i \Rightarrow \hat{\lambda}_{MLE}|_{\bar{x}} = \frac{1}{100} \cdot [0(143) + 1(156) + 2(68) + 3(27) + 4(5)] = 3.98 / 100 = 0.995$  (Ans)

Now Probability that no of RBC per cell is at most one is  $P(X=0) + P(X=1)$ . Now for this we will use the MLE so that we get MLE of the required probability

i.e. using  $\hat{\lambda}_{MLE}|_{\bar{x}}$  we get  $P(X=0) + P(X=1)$

$$= e^{-\hat{\lambda}_{MLE}} + \hat{\lambda}_{MLE} \cdot e^{-\hat{\lambda}_{MLE}} = e^{-\hat{\lambda}} (1 + \hat{\lambda}) = 0.737$$
 (Ans)

# Practical Day 8

pioneerpaper.co  
Page: Stat Proc  
Date: 19th April 2023

1) Below are given the yields in gm per plot (plot size = 1/2000 acre) for three varieties of seed cotton:

Variety 1	V2	V3
77	109	46
70	106	70
63	137	73
81	79	65
95	134	61
81	79	46
101	126	98

Test if the varieties differ significantly among themselves. Take  $\alpha = 0.05$

Anova 1-way layout. We have data on yields classified according to 3 varieties of seed cotton.

Hence we will use 1-way layout of ANOVA to model the data.

Total ~~21~~ plots have been observed where each variety has been yielded over 7 plots.

Let us consider the model,

$$y_{ij} = \mu + \alpha_i + \epsilon_j \quad j = 1, 2, \dots, n; \\ i = 1, 2, \dots, k$$

where  $y_{ij}$  : yields in  $j$ th plot receiving  $i$ th variety

$\mu$  : general effect.

$\alpha_i$  : additional effect due to  $i$ th variety  
 $i = 1, 2, \dots, k$

$\epsilon_j$  : random error associated with  $y_{ij}$

Assumptions  $\epsilon_j \sim N(0, \sigma^2)$

$K$  populations are homoscedastic

We are interested to see whether yields vary over ~~3~~ varieties of seed cotton. Hence, set of hypotheses will be

$$H_0: \alpha_1 = \alpha_2 = \dots = \alpha_k = 0 \quad \text{vs} \quad H_1: \text{not } H_0$$

Test statistic :  $F = \frac{MSB}{MSE} \sim F_{k-1, n-k}$

where,  $MSB = \sum_{i=1}^{k-1} n_i (\bar{y}_{i0} - \bar{y}_{00})^2$  and

$$MSE = \frac{\sum_{ij} (y_{ij} - \bar{y}_{i0})^2}{n-k}$$

A right tailed test based on  $F$  would be appropriate.  
We reject  $H_0$  at level  $\alpha$  if (observed  $F$ )  $> F_{\alpha; k-1; n-k}$   
where,

$F_{\alpha; k-1; n-k}$  denotes the upper  $\alpha$ -point of  
 $F_{k-1; n-k}$  distribution.

[ Steps : Stat  $\rightarrow$  Anova  $\rightarrow$  1-way...  $\rightarrow$  Fill in the  
Anova required fields ]

### Calculations

$K$  : no of varieties of seeds

$n_i$  : total no of plots of  $i^{th}$  variety seed  $i=1(1)3$

$\therefore n = \sum n_i = 21 >$  total no of plots.

$$\therefore F_{\alpha; k-1; n-k} = F_{0.05; 2; 19}$$

Using minitab we found the values of the following table

### Anova table

Source	df	Sum of squares	MS	Obsd $F$	$F_{k-1, n-k}$	Decision
Variety	2	SSB = 7089	MSB = 3544	9.85	3.55	Ref
Error	18	SSE = 6477	MSE = 360			
Total	20	TSS = 13566				

We have rejected  $H_0$  at 5% level of significance. Therefore,  
on the basis of the data, we can conclude  
at 5% level that yield varies significantly over  
varieties of seed cotton.

$$F_{\alpha; m;n} \Rightarrow P(F > F_{\alpha; m;n}) = \alpha \Rightarrow F_{\alpha; m;n} = P_{F_{m;n}}^{-1}(1-\alpha)$$

Similarly for t. as well and for any distn