

NOVEMBER 2020: IN SEMESTER ASSESSMENT, B.TECH, III-SEMESTER

TEST – 2

**UE19CS203 – STATISTICS FOR DATA SCIENCE
SCHEME & SOLUTION**

Time: 1.30 Hrs		Answer All Questions	Max Marks: 40
Q. No.			Marks
1.	a)	<p>A statistics instructor believes that fewer than sixty percentages of her university students attended the opening ceremony of a university fest. She surveys 64 of her students and finds that 36 attended the fest. Write appropriate null and alternative hypothesis.</p> <p>Solution:</p> <p>Null Hypothesis : $p \geq 0.6$ (Each hypothesis 1 mark)</p> <p>Alternate Hypothesis : $p < 0.6$</p>	2
	b)	<p>Globally the long-term proportion of newborns who are female is 51.46%. A researcher believes that the proportion of girls at birth changes under severe economic conditions. To test this belief randomly selected birth records of 5,000 babies born during a period of economic recession were examined. It was found in the sample that 52.55% of the newborns were girls. Determine whether there is sufficient evidence to support the researcher's belief.</p> <p>Solution:</p> <p>$H_0 : p = 0.5146$</p> <p>$H_1 : p \neq 0.5146$ (1 mark)</p> <p>Let $p_0 = 0.5146$</p> <p>$q_0 = 1 - p_0 = 1 - 0.5146 = 0.4854$</p> <p>$\hat{P} = 0.5255$</p> <p>Z statistic = $(\hat{p} - p_0) / \sqrt{p_0 * q_0 / n}$</p> <p>$= (0.5255 - 0.5146) / \sqrt{0.5146 * 0.4854 / 5000}$ (1 mark)</p> <p>$= 1.5352 \text{ or } 1.54$</p> <p>Since it is a two tailed test, $p\text{-value} = 2 * P(Z < -1.5352) = 2 * 0.0618 = 0.1236$ (1 mark)</p> <p>Since $p\text{-value} = 0.1236 \leq 0.05$ is false, we Don't reject H_0. (1 mark)</p>	4
	c)	<p>To compare customer satisfaction levels of two competing service providers of a particular service, 174 customers of Service Provider1 and 355 customers of Service Provider2 were randomly selected and were asked to rate their service providers on a five-point scale, with 1 being least satisfied and 5 most satisfied. Average rating obtained by service provider 1 was 3.51 with a standard deviation of 0.51 from their customers and the average rating obtained by service provider 2 was 3.24 with a standard deviation of 0.52. Test at the 5% level of significance whether the data provide sufficient evidence to conclude that Service provider1 has a higher mean satisfaction rating than does Service provider 2.</p> <p>Solution:</p>	4

		<p>Given $\bar{x} = 3.51$; $s_x = 0.51$ and $\bar{y} = 3.24$; $s_y = 0.52$ and $n_x = 174$; $n_y = 355$ Let μ_X = Mean service rating of Service Provider 1 Let μ_Y = Mean service rating of Service Provider 2</p> <p>$H_0 : \mu_X - \mu_Y \leq 0$ $H_1 : \mu_X - \mu_Y > 0$ (1 mark) Let $\mu_X - \mu_Y = 0$</p> <p>Z statistic = $(\bar{x} - \bar{y}) - (\mu_X - \mu_Y) / \sqrt{s_x^2/n_x + s_y^2/n_y}$ $= (3.51 - 3.24) - 0 / \sqrt{0.51^2/174 + 0.52^2/355} = 5.684$ (1 mark)</p> <p>Since it is a right tailed test, p-value = $P(Z > 5.684) = 1 - P(Z \leq 5.684)$ $= 1 - 1 = 0$ (1 mark)</p> <p>Since p-value = $0 \leq 0.05$, we reject H_0. (1 mark)</p>																																					
2.	a)	<p>A man who got transferred to a new place is trying to determine which of the two routes to work has the shorter average driving time. Times in minutes for six trips on route A and five trips on route B are follows: A: 16.0 15.7 16.4 15.9 16.2 16.3 B: 17.2 16.9 16.1 19.8 16.7 Can you conclude that the mean time is less for route A?</p> <p>μ_X – represents mean time for route B. μ_Y – represents mean time for route A. $H_0 : \mu_X \leq \mu_Y$ $H_1 : \mu_X > \mu_Y$ or 1 mark</p> <p>$H_0 : \mu_X - \mu_Y \leq 0$ $H_1 : \mu_X - \mu_Y > 0$</p> <table><thead><tr><th>Value</th><th>Rank</th><th>Sample</th></tr></thead><tbody><tr><td>15.7</td><td>1</td><td>Y</td></tr><tr><td>15.9</td><td>2</td><td>Y</td></tr><tr><td>16</td><td>3</td><td>Y</td></tr><tr><td>16.1</td><td>4</td><td>X</td></tr><tr><td>16.2</td><td>5</td><td>Y</td></tr><tr><td>16.3</td><td>6</td><td>Y</td></tr><tr><td>16.4</td><td>7</td><td>Y</td></tr><tr><td>16.7</td><td>8</td><td>X</td></tr><tr><td>16.9</td><td>9</td><td>X</td></tr><tr><td>17.2</td><td>10</td><td>X</td></tr><tr><td>19.8</td><td>11</td><td>X</td></tr></tbody></table> <p>1 mark</p> <p>Here, $W = 42$ (Counting all X(route B) ranks) 1 mark Since $m = 5$ and $n = 6$, we find that the area to the left of $W = 42$ is 0.0152 1 mark Since $P < 0.05$, we reject H_0 and conclude that mean lifetime for route A is less. 1 mark</p>	Value	Rank	Sample	15.7	1	Y	15.9	2	Y	16	3	Y	16.1	4	X	16.2	5	Y	16.3	6	Y	16.4	7	Y	16.7	8	X	16.9	9	X	17.2	10	X	19.8	11	X	5
Value	Rank	Sample																																					
15.7	1	Y																																					
15.9	2	Y																																					
16	3	Y																																					
16.1	4	X																																					
16.2	5	Y																																					
16.3	6	Y																																					
16.4	7	Y																																					
16.7	8	X																																					
16.9	9	X																																					
17.2	10	X																																					
19.8	11	X																																					

b)	Given the following contingency table for hair colour and eye colour. Find the value of Chi-Square and is there any good association between the two.	5																																							
	<table><tr><td>Hair colour \ Eye colour</td><td>Fair</td><td>Brown</td><td>Black</td></tr><tr><td>Grey</td><td>20</td><td>10</td><td>20</td></tr><tr><td>Brown</td><td>25</td><td>15</td><td>20</td></tr><tr><td>Black</td><td>15</td><td>5</td><td>20</td></tr></table> <p>Solution:</p> <table><tr><td>Hair colour \ Eye colour</td><td>Fair</td><td>Brown</td><td>Black</td><td>Row Sum</td></tr><tr><td>Grey</td><td>20(Observed) 20(Expected)</td><td>10(Observed) 10(Expected)</td><td>20(Observed) 20(Expected)</td><td>50</td></tr><tr><td>Brown</td><td>25 24(Expected)</td><td>15 12(Expected)</td><td>20 24(Expected)</td><td>60</td></tr><tr><td>Black</td><td>15 16(Expected)</td><td>5 8(Expected)</td><td>20 16(Expected)</td><td>40</td></tr><tr><td>Column Sum</td><td>60</td><td>30</td><td>60</td><td>150</td></tr></table> <p>(Expected values – 2 marks) Dof = (r-1)*(c-1) = (3-1)*(3-1)=2*2=4 0.5 mark Calculated χ^2 value = $\sum \sum (\text{Observed} - \text{Expected})^2 / \text{Expected}$ $= (20-20)^2/20 + (10-10)^2/10 + (20-20)^2/20 + (25-24)^2/24 + (15-12)^2/12 + (20-24)^2/24 + (15-16)^2/16 + (5-8)^2/8 + (20-16)^2/16$ $= 1/24 + 9/12 + 16/24 + 1/16 + 9/8 + 1 = 3.6458$ 1.5 marks</p> <p>Calculated χ^2 value = 3.6458 Tabulated Value = 9.488 0.5 mark (at 5% level of significance with 4 degrees of freedom) Since, Calculated value < Tabulated value, Don't Reject Ho (Null hypothesis) 0.5 mark</p>	Hair colour \ Eye colour	Fair	Brown	Black	Grey	20	10	20	Brown	25	15	20	Black	15	5	20	Hair colour \ Eye colour	Fair	Brown	Black	Row Sum	Grey	20(Observed) 20(Expected)	10(Observed) 10(Expected)	20(Observed) 20(Expected)	50	Brown	25 24(Expected)	15 12(Expected)	20 24(Expected)	60	Black	15 16(Expected)	5 8(Expected)	20 16(Expected)	40	Column Sum	60	30	60
Hair colour \ Eye colour	Fair	Brown	Black																																						
Grey	20	10	20																																						
Brown	25	15	20																																						
Black	15	5	20																																						
Hair colour \ Eye colour	Fair	Brown	Black	Row Sum																																					
Grey	20(Observed) 20(Expected)	10(Observed) 10(Expected)	20(Observed) 20(Expected)	50																																					
Brown	25 24(Expected)	15 12(Expected)	20 24(Expected)	60																																					
Black	15 16(Expected)	5 8(Expected)	20 16(Expected)	40																																					
Column Sum	60	30	60	150																																					
3.	a) Given the null hypothesis: that a process is producing no more than the maximum allowable rate of defective items. Write the Type II Error(as a statement).	2																																							
	<p>Solution:</p> <p>Ho = process producing no more than k defectives Ho false means process producing more than k defectives.</p> <p>Type II error = P(do not reject Ho when Ho is false)</p> <p>Type II error: The process is not producing too many defectives when it actually is.</p>																																								

	b)	<p>A copper smelting process is supposed to reduce the arsenic content of the copper to less than 1000 ppm. Let μ denote the mean arsenic content for copper treated by this process, and assume that the standard deviation of arsenic content is $\sigma = 100$ ppm. The sample mean arsenic content X of 75 copper specimens will be computed, and the null hypothesis $H_0 : \mu \geq 1000$ will be tested against the alternate $H_1 : \mu < 1000$.</p> <p>i) A decision is made to reject H_0 if $X \leq 980$. Find the level of this test.</p> <p>ii) Find the power of the test in part (i) if the true mean content is 965 ppm.</p> <p>Solution:</p> <p>Given n= sample size = 75 True mean = 965 $\sigma = 100$; $\alpha = 5\% = 0.05$ $H_0: : \mu \geq 1000$ $H_1: \mu < 1000$</p> <p>i) $z = (\bar{x} - \mu) / \sigma / \sqrt{n} = (980 - 1000) / (100 / \sqrt{75}) \approx -1.73$ The level of the test is the probability of rejecting the null hypothesis. Level = $P(Z < -1.73) = 0.0418$ 2 marks</p> <p>ii) The power is the probability of rejecting the null hypothesis when the alternative hypothesis is true $z = (\bar{x} - \mu) / (\sigma / \sqrt{n}) = (980 - 965) / (100 / \sqrt{75}) \approx 1.30$ Power= $P(Z < 1.30)=0.9032$ 2 marks</p>	4																																										
	c)	<p>A certain type of seed has always grown to a mean height of 9.5 inches, with a standard deviation of 1 inch. Based on past experiment, the mean height of a seed is known to be distributed approximately normal. A researcher wishes to find out whether some new enriched conditions would improve the mean height. He wants to use $\alpha = .01$ test and would like to have a 96% chance of rejecting the null hypothesis if the mean height is 10.5 inches. Determine the sample size for the test.</p> <p>Solution:</p> <p>Given $\alpha = .01$; $1 - \beta = 0.96$ i.e. $\beta = 0.04$ 1 mark</p> <p>$\mu_{H_0} = 9.5$ and $\mu_{H_A} = 10.5$ 1 mark</p> <p>We know that $n = (\sigma^2 * (Z_\alpha + Z_\beta)^2) / (\mu_{H_A} - \mu_{H_0})^2$ 1 mark $n = 12 * (2.33 + 1.75)^2 / (10.5 - 9.5)^2 = 16.65 \approx 17$ $n = 17$ 1 mark</p>	4																																										
4.	a)	<p>Find the least-squares regression line for the given data:</p> <table><tr><td>x</td><td>2</td><td>2</td><td>6</td><td>8</td><td>10</td></tr><tr><td>y</td><td>0</td><td>1</td><td>2</td><td>3</td><td>3</td></tr></table> <p>Solution:</p> <table><tr><td>x</td><td>y</td><td>$x - \bar{x}$</td><td>$y - \bar{y}$</td><td>$(x - \bar{x})^2$</td><td>$(x - \bar{x}) * (y - \bar{y})$</td></tr><tr><td>2</td><td>0</td><td>-3.6</td><td>-1.8</td><td>12.96</td><td>6.48</td></tr><tr><td>2</td><td>1</td><td>-3.6</td><td>-0.8</td><td>12.96</td><td>2.88</td></tr><tr><td>6</td><td>2</td><td>0.4</td><td>0.2</td><td>0.16</td><td>0.08</td></tr><tr><td>8</td><td>3</td><td>2.4</td><td>1.2</td><td>5.76</td><td>2.88</td></tr></table>	x	2	2	6	8	10	y	0	1	2	3	3	x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(x - \bar{x}) * (y - \bar{y})$	2	0	-3.6	-1.8	12.96	6.48	2	1	-3.6	-0.8	12.96	2.88	6	2	0.4	0.2	0.16	0.08	8	3	2.4	1.2	5.76	2.88	5
x	2	2	6	8	10																																								
y	0	1	2	3	3																																								
x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(x - \bar{x}) * (y - \bar{y})$																																								
2	0	-3.6	-1.8	12.96	6.48																																								
2	1	-3.6	-0.8	12.96	2.88																																								
6	2	0.4	0.2	0.16	0.08																																								
8	3	2.4	1.2	5.76	2.88																																								

		10	3	4.4	1.2	19.36	5.28	
		<p align="center">(Table 1 mark)</p> <p>$\bar{x} = 5.6$ and $\bar{y} = 1.8$ (0.5 mark)</p> <p>$\sum(x-\bar{x})^2 = 51.2$; (0.5 mark)</p> <p>$\sum(x-\bar{x})(y-\bar{y}) = 17.6$ (0.5 mark)</p> <p>$\hat{\beta}_1 = \sum(x-\bar{x})(y-\bar{y}) / \sum(x-\bar{x})^2 = 17.6 / 51.2 = 0.34375$ (1 mark)</p> <p>$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 1.8 - 0.34375 * 5.6 = -0.125$ (1 mark)</p> <p>Least Squares Regression Line is $\hat{y} = 0.34375 x - 0.125$ (0.5 mark)</p>						
b)	Find the uncertainties in the estimates of $\hat{\beta}_0$ and $\hat{\beta}_1$ of part a).	<p>$\sum(y-\bar{y})^2 = (-1.8)^2 + (-0.8)^2 + (0.2)^2 + (1.2)^2 + (1.2)^2 = 6.8$ (1 mark)</p> <p>Correlation Coefficient $r = \sum(x-\bar{x})(y-\bar{y}) / (\sqrt{\sum(x-\bar{x})^2} * \sqrt{\sum(y-\bar{y})^2})$</p> <p>$r = 17.6 / \sqrt{51.2} * \sqrt{6.8} = 0.9432$ (1 mark)</p> <p>Error Uncertainty $s = \sqrt{((1-r^2) * (\sum(y-\bar{y})^2) / (n-2))}$ (1 mark)</p> <p>$n=5$</p> <p>$s_{\hat{\beta}_0} = s * \sqrt{(1/n) + (\bar{x}^2 / \sum(x-\bar{x})^2)}$ (1 mark)</p> <p>$s_{\hat{\beta}_1} = s / \sqrt{\sum(x-\bar{x})^2} = 0.5002 / \sqrt{51.2} = 0.0699$ (1 mark)</p>						5