

STA501. Flowchart with formulae

1		Is the data?	
	Quantitative	(NOTE: If σ given, must use z , not t)	go to 2
	Qualitative		go to 18
2	Quantitative	How many samples are there ?	
		Population mean – One	goto 3
		Population mean – Two	goto 8
		Population mean – Many	goto 14
		Relationship between two quantitative variables	goto 24
3	Quantitative One Sample	Is $n \geq 30$?	
		Yes	goto 4
		No	goto 5
4	Quantitative One Sample $N \geq 30$	Confidence Interval $\bar{x} \pm Z_{\alpha/2} \times \frac{s}{\sqrt{n}}$ For interpretation of Confidence Interval (CI)	goto 29
		Hypothesis Test $Z = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \sim N(0,1)$ $X \sim N(\mu, \sigma^2)$	

STA501. Flowchart with formulae

5	Quantitative One Sample N < 30	Is data normally distributed? Yes No	(Use QQ Plot) goto 7 goto 6
6	Quantitative One Sample N < 30 Data not normally distributed	Assumption not met.	
		state that assumptions not met and therefore results not reliable /accurate)	goto 7
7	Quantitative One Sample N < 30 Data normally distributed	Confidence Interval $x \pm t_{n-1, \frac{\alpha}{2}} \times \frac{s}{\sqrt{n}}$ Hypothesis Test $t = \frac{x - \mu}{\frac{s}{\sqrt{n}}} \sim t_{n-1}$	

STA501. Flowchart with formulae

8	Quantitative Two samples	Is data paired or independent?
		<p>Paired</p> <p style="padding-left: 40px;">Calculate difference of means, and treat as single sample goto 3</p> <p><i>Note : paired data : data is related / dependent Different sample size indicates independent samples.</i></p>
		Independent goto 9
9	Quantitative Two samples Independent	Are both samples ≥ 30
		Yes goto 10
		No goto 11
10	Quantitative Two samples Independent Both samples ≥ 30	<p>Confidence Interval</p> $(\bar{x}_1 - \bar{x}_2) \pm z_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$ <p>Hypothesis Test</p> $Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$

STA501. Flowchart with formulae

11	Quantitative Two samples Independent One / Both Samples n < 30	<p>Test if Variances are equal</p> $H_0: \sigma_1^2 = \sigma_2^2$ $H_1: \sigma_1^2 \neq \sigma_2^2$ $F = \frac{S_1^2}{S_2^2} \sim F_{n_1-1, n_2-1} \text{ (biggest on top, smallest on bottom)}$ <p>Reject H_0 if $F_{obs} > F_{n_1-1, n_2-1, \alpha/2}$</p> <p>Are variances equal?</p>
		Yes goto 12
		No goto 13
12	Quantitative Two samples Independent One / Both samples n < 30 Variances equal	<p>Confidence Interval</p> $(\bar{x}_1 - \bar{x}_2) \pm t_{df, \alpha/2} \sqrt{s_p^2 \left(\frac{1}{n} + \frac{1}{n} \right)}$ $df = (n_1 + n_2 - 2)$ <p>Hypothesis Test</p> $t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{s_p^2 \left(\frac{1}{n} + \frac{1}{n} \right)}} \sim t_{n_1+n_2-2}$ $s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$
13	Quantitative Two samples Independent One / Both Samples n < 30 Variances NOT equal	<p>Confidence Interval</p> $(\bar{x}_1 - \bar{x}_2) \pm z_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$ $df = \min(n_1 - 1, n_2 - 1)$ <p>Hypothesis Test</p> $t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{s_p^2 \left(\frac{1}{n} + \frac{1}{n} \right)}} \sim t_{\min(df)}$

STA501. Flowchart with formulae

14	Quantitative Many Samples	<p>ANOVA Test</p> <p>Test if ANOVA Test assumptions met:-</p> <p>Independent Simple Random Sample (SRS)</p> <p>Parent populations are normally distributed use QQ Plots / boxplots</p> <p>All populations have same standard deviation practical rule</p> <p>Results are approximately correct if largest sample std dev is no more than twice as large as smallest sample std dev (or side by side boxplots)</p> <p>Variability Between Samples</p> $F = \frac{MSG}{MSE} \sim F_{dfG, dfE}$ <p>Reject H0 if $F_{obs} \geq F_{dfG, dfE}$</p> <p>Hypothesis Testing</p> $H_0 : \mu_1 = \mu_2 = \mu_3$ $H_1 : \text{not all } \mu' \text{ sare equal}$ <p>Is the NULL Hypothesis Rejected?</p>
		Yes goto 15
		No <i>no further test required</i>
15	Quantitative Many Samples not all μ 's equal	<p>Use Bonferroni Method</p> <p>valid only if Ho rejected</p> <p>not all μ's are equal</p> <p>Post-hoc procedure.</p>

STA501. Flowchart with formulae

18	Qualitative	Is there one or two samples?
		Population proportion - One sample goto 19
		Population proportion - Two samples goto 21
		Relationship between 2 categorical variables goto 23a
		One qualitative variable with many levels – Goodness of Fit goto 23b
19	Qualitative One sample	$np > 5$ $n(1 - p) > 5$
		Yes goto 20
		No got 20, <i>NOTE: results not reliable</i>
20	Qualitative One sample np and nq > 5	<p>Confidence Interval</p> $\hat{p} \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$ <p>Hypothesis Test</p> $Z = \frac{\hat{p} - p}{\sqrt{\frac{p(1 - p)}{n}}} \sim N(0,1)$
21	Qualitative Two samples	<p>Both Samples</p> $np > 5$ $n(1 - p) > 5$
		Yes goto 22
		No goto 22. <i>NOTE results not reliable</i>
22	Qualitative Two samples np and nq > 5	<p>Confidence Interval</p> $(\hat{p}_1 - \hat{p}_2) \pm Z^* \sqrt{\hat{p}(1 - \hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$ <p>\hat{p} = count of success in both samples / total count in both samples</p> <p>Hypothesis Test</p> $Z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\hat{p}(1 - \hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$

23a	Chi square test of independence	<p> H_0 : variables are NOT related and are independent H_1 : variables are related and dependent </p> $X^2 = \sum \frac{(\text{Observed Value})^2}{\text{Expected Value}} - N \sim X_{df}^2$ <p> $df = (\text{rows} - 1) \times (\text{columns} - 1)$ </p> <p> <i>Reject H_0 if $X_{obs}^2 > X_{df,\alpha}^2$</i> </p> <p>Steps For Chi-Square Test</p> <p>Table 1 : Observed Values Table 2 : Expected Values (E)</p> $\text{Expected Value} = \frac{(\text{Row Total}) \times (\text{Column Total})}{\text{Grand Total}}$ $\text{Expected Value} = \frac{\text{Row Total}}{\text{Grand Total}} \times \frac{\text{Column Total}}{\text{Grand Total}} \times \text{Grand Total}$ <p>Assumptions</p> <ul style="list-style-type: none"> • No table cell has an Expected Value less than 1 • No more than 20% of cells have expected values less than 5 • Outcome of each trial can only fall into one of the categories
23b	Chi square test goodness of fit	<p>Does data match a pattern?</p> <p> H_0 : observed value DOES match expected value H_1 : observed value DOES NOT match expected value </p> $X^2 = \sum \frac{(\text{Observed Value})^2}{\text{Expected Value}} - N \sim X_{df}^2$ <p> $df = n - \text{number of assumptions}$ $df = n - 1$ </p> <p> <i>Reject H_0 if $X_{obs}^2 > X_{df,\alpha}^2$</i> </p> <p>Steps For Chi-Square Test</p> <p>Table 1 : Observed Values Table 2 : Expected Values (E) apply pattern to N to determine this???</p> <p>Assumptions</p> <ul style="list-style-type: none"> • No table cell has an Expected Value less than 1 • No more than 20% of cells have expected values less than 5 • Outcome of each trial can only fall into one of the categories

STA501. Flowchart with formulae

24	Correlation	goto 25
	Regression	goto 26
25	Correlation	<p>$H_0 : \rho = 0$ (linear relationship – does not exist)</p> <p>$H_1 : \rho \neq 0$ (linear relationship – does exist)</p> $t = \frac{r - \rho}{\sqrt{\frac{1 - r^2}{n - 2}}} \sim t_{n-2}$ <p>$df = n - 2$ ($n = \text{total number of points} \mid \text{pairs} \mid y \text{ values}$)</p> <p>Coefficient of Determination (CoD)</p> $R = r^2$ <p>Coefficient of Determination (CoD) tells us that</p> <ul style="list-style-type: none"> • x% of the variation in y is • explained by the linear regression on x.
26	Regression	$\hat{y} = a + bx$ <ul style="list-style-type: none"> • b : the slope of the estimated regression line <ul style="list-style-type: none"> ○ The addition of one unit of x is associated with a change of b units of y on average. • A : the intercept of the estimated regression line <ul style="list-style-type: none"> ○ When x is zero, y is a. /gives the minimum y value. <p>Using regression equation for estimation – should be within range of existing data.</p> <p>Confidence Interval</p> $b \pm t_{n-2, \alpha/2} SE_b$ <p>Hypothesis Test</p> <p>$H_0 : \beta = 0$ (slope is zero)</p> <p>$H_1 : \beta \neq 0$ (slope is NOT zero) (linear relationship between the two variables)</p> $t_{obs} = \frac{b}{SE_b} \sim t_{n-2}$ <p>SEb - does not need to be calculated This value will always be given (prob as part of S-Plus output)</p> <p>$tolerance = 1 - r^2$ $tolerance < 0.1$, suggests variables are colinear</p> <p>Used for determining best model</p> $Adjusted R^2 = 1 - (1 - Multiple R^2) \frac{n - 1}{n - p - 1}$

STA501. Flowchart with formulae

27	Hypothesis Template	<ol style="list-style-type: none"> 1. Write null and alternative hypothesis <ul style="list-style-type: none"> ○ H_0: does NOT have any affect ○ H_1 : DOES have an affect 2. State level of significance <ul style="list-style-type: none"> ○ α = Level of significance ○ $CI = (100 - \alpha)$ Confidence Interval 3. State the sampling distribution / test statistic <ul style="list-style-type: none"> ○ Write down Actual Formula / Equation ○ Write down Name of Formula / Equation 4. Decision Rule: Find Critical region (use diagram) <i>Examples</i> <ul style="list-style-type: none"> ○ t_{cv} OR ○ Reject H_0 if p-value < alpha 5. Calculate <ul style="list-style-type: none"> ○ Test statistic (<i>results from using formula</i>) ○ Observed value 6. Decision (H_0 : <i>Reject or Accept</i>) <ul style="list-style-type: none"> ○ Explicitly compare observed value with critical value OR ○ Explicitly compare <i>p-value</i> and <i>alpha</i> 7. Conclusion <ul style="list-style-type: none"> ○ Interpretation of test result ○ Answer in terms of original problem. 														
28	ANOVA	<ol style="list-style-type: none"> 1 For each sample <table style="width: 100%; border: none;"> <tr> <td style="width: 50%;">sample size n</td> <td style="width: 50%;">also N</td> </tr> <tr> <td>Total T</td> <td>also Grand Total</td> </tr> <tr> <td>Total²</td> <td>T²</td> </tr> <tr> <td>T² / n</td> <td>also Total</td> </tr> </table> 2 Calculated sum of al X² (ΣX^2) 3 <table style="width: 100%; border: none;"> <tr> <td style="width: 50%;">$SST = (\Sigma X^2)$</td> <td style="width: 50%;">$(\text{Grand Total})^2 / N$</td> </tr> <tr> <td>$SSG = \Sigma T^2/n$</td> <td>$(\text{Grand Total})^2 / N$</td> </tr> <tr> <td>$SSE = SST - SSG$</td> <td></td> </tr> </table> 4 <p><i>df treatment = number of samples - 1</i> <i>df total = N - 1</i> <i>df error = df total - df treatment</i></p> 	sample size n	also N	Total T	also Grand Total	Total ²	T ²	T ² / n	also Total	$SST = (\Sigma X^2)$	$(\text{Grand Total})^2 / N$	$SSG = \Sigma T^2/n$	$(\text{Grand Total})^2 / N$	$SSE = SST - SSG$	
sample size n	also N															
Total T	also Grand Total															
Total ²	T ²															
T ² / n	also Total															
$SST = (\Sigma X^2)$	$(\text{Grand Total})^2 / N$															
$SSG = \Sigma T^2/n$	$(\text{Grand Total})^2 / N$															
$SSE = SST - SSG$																
29	CI	The true population mean will be contained in the xx% confidence interval, xx% of the time.														