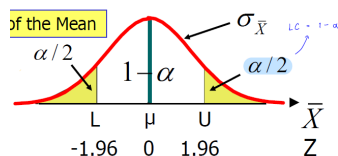


Confidence Interval (CI)

- Assuming data is normally distributed or if not use $n > 30$
- CI of 95% is that 95% of the population is in this CI meaning the population parameters (μ) also lies within this interval.
- Level of confidence (α), $\frac{\alpha}{2}$ for the proportion in the upper and lower tail
 - o 95% CI $\rightarrow \alpha/2 = 0.025 \rightarrow Z = \pm 1.96$
 - o 99% CI $\rightarrow \alpha/2 = 0.05 \rightarrow Z = \pm 2.58$
- For one side (Lower CI or Upper CI), just use α
- o 95% Upper CI $\rightarrow \alpha = 0.05 \rightarrow Z = 1.645$



CI - known σ

- Use Z-test.
- $\mu \leq \bar{x} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \rightarrow \mu \in [\bar{x} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \bar{x} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}]$
- Lower confidence bound: $\mu \geq \bar{x} - Z_{\alpha} \cdot \frac{\sigma}{\sqrt{n}} \rightarrow [\mu, \infty)$
- Upper confidence bound: $\mu \leq \bar{x} + Z_{\alpha} \cdot \frac{\sigma}{\sqrt{n}} \rightarrow (-\infty, \mu]$
- $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$
- Determine sample size of mean. $n = \frac{Z_{\alpha/2}^2 \cdot \sigma^2}{\text{error}^2}$

CI - unknown σ & small n ($n < 30$)

- Use T-test. $t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$, with $DOF = (n - 1)$
- $\mu \leq \bar{x} \pm t_{\alpha/2, n-1} \cdot \frac{s}{\sqrt{n}}$
- T-distribution approach normal distribution when n is large.

CI - unknown σ & large n ($n \geq 30$)

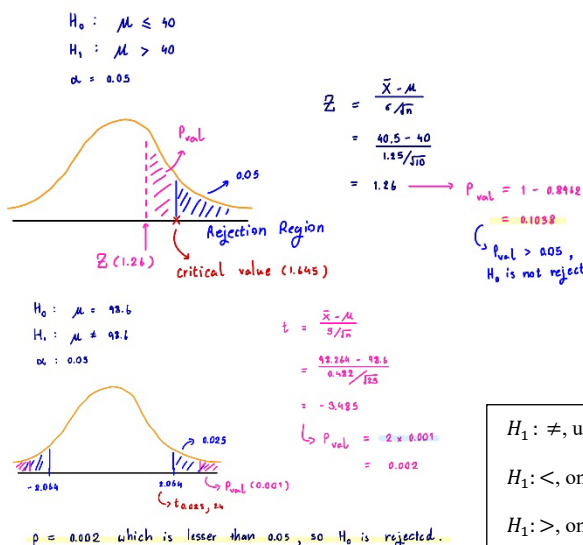
- As n is large, we can assume that the population is normally distributed, so the Z-test is used. $Z = \frac{\bar{x} - \mu}{s/\sqrt{n}} \rightarrow \text{noted: when } n \text{ is large, } \sigma \approx S = \sqrt{\frac{\sum(x_i - \bar{x})^2}{n-1}}$
- $\mu \leq \bar{x} \pm Z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}$

CI - population proportion

- Binary outcome (binomial distribution)
- Assumed normal distributed when $np \geq 5, n(1-p) \geq 5$, Z-test is used.
- $p \leq \bar{p}_s \pm Z_{\alpha/2} \cdot \sqrt{\frac{\bar{p}_s(1-\bar{p}_s)}{n}} \rightarrow p_s = \frac{n \text{ of success}}{n \text{ total sample}}$
- Determine sample size. $n = \frac{Z_{\alpha/2}^2 \cdot p_s(1-p_s)}{\text{error}^2}$

Hypothesis Testing (1 sample)

- Idea is to try proving the null hypothesis (H_0) using a given sample.
- Test method.
 - o Z-test:
 - known $\sigma \rightarrow Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$
 - proportion (p_s) $\rightarrow Z = \frac{p_s - p}{\sqrt{p(1-p)/n}} \rightarrow \text{Noted: } \sigma_{p_s} = \sqrt{\frac{p(1-p)}{n}}, \mu_{p_s} = p$
 - o T-test: unknown $\sigma \rightarrow t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$
- If the value from test method (Z, t) is not in the rejection zone created by the level of significance (α), H_0 is accepted or known as "There is not enough evidence to reject H_0 ".
- Critical value: value from level of significance (α) at the edge of rejection zone (\sim CI).
- Null hypothesis (H_0) is the hypothesis that contains "=". Ex. $\mu = 7, \mu \geq 9, \mu \leq 12$.
- Alternate hypothesis (H_1) is the opposite of null hypothesis.
- Steps
 1. Set H_0
 2. Cal for test statistics (Z, t, f)
 3. Determine region of rejection & critical value from α
 4. Compare value.
 - a. Test stat: If test stat lies in region of rejection (RoR), H_0 is rejected.
 - b. P-value: convert test stat into probability, if probability is less than α , H_0 is rejected. For two side, the P-value is doubled.

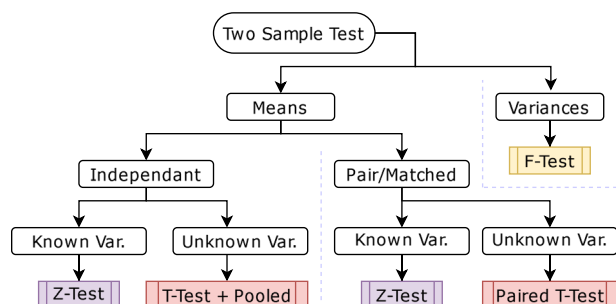


$H_1: \neq$, use two tails.

$H_1: <$, one tail, RoR on left.

$H_1: >$, one tail, RoR on right.

Hypothesis Testing (2 sample)



- **Criteria:** Assumed population is normally distributed or $n > 30$

Mean - Independent Sample

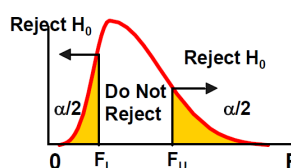
- Variance known (Z-test) $\rightarrow Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$
- Variance unknown (Pooled T-test)
 - o Population is normal or $n > 30$.
 - o Assumed equal ($\sigma_1 = \sigma_2$)
 - o Pooled sample variance $\rightarrow S_p^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{(n_1-1) + (n_2-1)}$
 - o $t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{S_p^2(\frac{1}{n_1} + \frac{1}{n_2})}}$, using $DOF = n_1 + n_2 - 2$

Mean - Related Sample (Paired/Match)

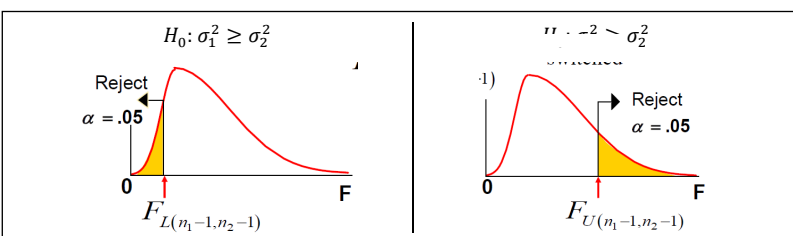
- Must observed paired (ex. Before & after)
- $n_1 = n_2 = n$ always
- Use the difference between pairs.
 - o $D_i = x_{1,i} - x_{2,i}$
 - o $\bar{D} = \frac{\sum D_i}{n}$
- Variance known $\rightarrow Z = \frac{\bar{D} - \mu_D}{\sigma_D/\sqrt{n}}$
- Variance unknown
 - o $t = \frac{\bar{D} - \mu_D}{s_D/\sqrt{n}}$, using $DOF = n - 1$
 - o $S_D = \sqrt{\frac{\sum(D_i - \bar{D})^2}{n-1}}$

Variance

- Both populations are normally distributed, this test is not robust to this violation.
- Use F-test: $F = \frac{S_1^2}{S_2^2}$
- $DOF_1 = n_1 - 1, DOF_2 = n_2 - 1$
- Critical value: $F_{L(DOF_1, DOF_2)} = \frac{1}{F_{U(DOF_2, DOF_1)}}, F_{U(DOF_1, DOF_2)}$
- If F - value is between F_L & F_U , H_0 is not rejected.



Note: The F-test graph always starts from the right.



ANOVA

- Control one or more independent variables: **treatment factors**.
- Observe effects as dependent variables: **response variable**.
- Condition
 - o Samples are randomly and independently drawn.
 - o Population is normally distributed.
 - Less sensitive when n is the same.
 - o Population variances are equal.
- $H_0: \mu_1 = \mu_2 = \dots = \mu_n$
 - o All μ is the same, **no treatment effect** → **the factor has no effect on the response variable**.
- H_1 : not all μ are the same.
 - o At least one μ is different, **there are treatment effect**.
- **Total variation** → $SST = \sum_{j=1}^c \sum_{i=1}^{n_j} (x_{ij} - \bar{x})^2$
 - $(x_{11} - \bar{x})^2 + (x_{12} - \bar{x})^2 + \dots + (x_{n_c c} - \bar{x})^2$
 - o x_{ij} : data i in group j
 - o n_j : number of observations in group j
 - o c: number of groups.
 - o \bar{x} : overall mean $\left(\frac{\text{sum of all } x_{ij}}{\text{total observation (n=c} \cdot \text{n}_j)} \right)$
 - o $SST = SSA + SSW$
 - o $DOF = n - 1$
- **Among-Group variation** → $SSA = \sum_{j=1}^c n_j (\bar{x}_j - \bar{x})^2$
 - $n_1(\bar{x}_1 - \bar{x})^2 + n_2(\bar{x}_2 - \bar{x})^2 + \dots + n_c(\bar{x}_c - \bar{x})^2$
 - o \bar{x}_j : sample mean of group j.
 - o $DOF = c - 1$
 - o $MSA = \frac{SSA}{c-1}$
- **Within-Group variation** → $SSW = \sum_{j=1}^c \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2$
 - $(x_{11} - \bar{x}_1)^2 + (x_{12} - \bar{x}_1)^2 + \dots + (x_{n_c c} - \bar{x}_c)^2$
 - o $DOF = n - c$
 - o $MSW = \frac{SSW}{n-c}$
- **F-test**: $f = \frac{MSA}{MSW}$, using $DOF_1 = c - 1, DOF_2 = n - c$
- Using f table if $f > F_{\alpha, DOF_1, DOF_2}$, reject H_0
- To find which μ is different, we use LSD.
 - o Compare pairwise, find absolute mean difference (AMD) of each pair.
 - Ex: $|\bar{x}_1 - \bar{x}_2|, |\bar{x}_1 - \bar{x}_3|, |\bar{x}_2 - \bar{x}_3|, \dots$
 - o Same sample size: $LSD = t_{\alpha/2, c(n_j-1)} \cdot \sqrt{\frac{2MSW}{n_j}}$
 - o Different sample size: $LSD = t_{\alpha/2, n-c} \cdot \sqrt{MSW \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}$
 - o If AMD of a pair is larger than LSD, that pair have different μ .

Two ways ANOVA

- Two factors affecting one response.

Regression

- Dependent variable (Y): response variable
 - Independent variable (X): explanatory or predictor variable
 - Idea: try to find a straight line that best fits the given data, **having smallest error** when calculated between the actual and predicted values.
 - **To determine the model** → $\hat{y} = b_0 + b_1x$
 - o \hat{y} : **predicted value (response)**
 - o \bar{y} : **average y** → $\bar{y} = \frac{\sum y_i}{n}$
 - o y : **actual value of y from dataset**
 - o x : **actual value of x from dataset (predictor)**
 - o b_1 : slope of the regression → $b_1 = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} = \frac{\sum xy - n\bar{x}\bar{y}}{\sum x^2 - n\bar{x}^2}$
 - o b_0 : y interception (AKA. bias) → $b_0 = \bar{y} - b_1\bar{x}$
 - Error (residual): $e = y - \hat{y}$
 - **To determine the fit of the model**, we calculate **Coefficient of determination (r^2)**
 - o $r^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$
 - o It explains how much the variability of **y** can be explained by the model based on **x**.
Ex. $r^2 = 0.69 \rightarrow 69\%$ of the data can be explained by x based model.
 - SST : Total variability about the mean → $SST = \sum (y - \bar{y})^2$
 - SSE : Variability about the regression line → $SSE = \sum (y - \hat{y})^2$
 - SSR : Total variability that explain the model → $SSR = \sum (\hat{y} - \bar{y})^2$
- $SST = SSR + SSE$
-
- Coefficient of correlation (r) is the expression of the strength of linear relationship.
 - o Value between -1 to $+1$
 - **To determine the significance of the model**, we use **F-test**.
 - $H_0: b_0 = 0 \rightarrow$ there is no relationship between x and y .
 - $MSR = \frac{SSR}{k}$, $MSE = \frac{SSE}{n-k-1} \rightarrow MSE$ is also known as σ^2
 - o k: number of independent variables
 - o n: number of observations
 - $f = \frac{MSR}{MSE}$, using $DOF_1 = k, DOF_2 = n - k - 1$
 - o Probability of f is known as **Significance F (F')**
 - o lower F' means the model is useful in predicting y.
 - o $F' < \alpha$, reject H_0 , there are relationship between x and y .

Multiple Regression analysis

- Study the model with multiple independent variables.
- $\hat{y} = b_0 + b_1x_1 + b_2x_2 + \dots + b_kx_k + \varepsilon$
 - o ε : random error
- The largest b_i show that x_i has **the highest impact** to \hat{y}
- r^2 explain how much \hat{y} can be explain by the model based on x_1, x_2, \dots, x_k
The rest $(1 - r^2)$ explain how the remaining \hat{y} can be explained by other variables that aren't included in the model.

	A	B	C	D	E	F	G	H	I
16									
17	SUMMARY OUTPUT								
18									
19	Regression Statistics								
20	Multiple R	0.81968							
21	R Square	0.67188							
22	Adjusted R Square	0.61222							
23	Standard Error	24312.6							
24	Observations	14							
25									
26	ANOVA								
27		df	SS	MS	F	Significance F			
28	Regression	2	13313936968	6.7E+09	11.2619				
29	Residual	11	6502131603	5.9E+08					
30	Total	13	19816068571						
31									
32		Coefficient	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
33	Intercept	146631	25482.0829	5.7543	0.0001	90545.2073	202716.5798	90545.2073	202716.5798
34	SF	43.8194	10.2810	4.2622	0.0013	21.1911	66.4476	21.1911	66.4476
35	AGE	-2898.69	796.5649	-3.6390	0.0039	-4651.9139	-1145.4586	-4651.9139	-1145.4586

The coefficient of determination (r^2) is 0.67.

A low significance level for F proves a relationship exists.

The regression coefficients are found here.

The P-values are used to test the individual variables for significance.

$$\hat{Y} = 146631 + 44X_1 - 2899X_2$$