

STA 106: Applied Statistical Methods: Analysis of Variance

Course Material Summary

University of California at Davis

Last Edit Date: 02/19/2022

Disclaimer and Term of Use:

1. We do not guarantee the accuracy and completeness of the summary content. Some of the course material may not be included, and some of the content in the summary may not be correct. You should use this file properly and legally. We are not responsible for any results from using this file.
2. This personal note is adapted from *Professor Xuehen Shi*. Please [contact us](#) to delete this file if you think your rights have been violated.
3. This work is licensed under a [Creative Commons Attribution 4.0 International License](#).

STA 106 Summary

Topic	Content
Statistical Concepts	<p>A population is the set of all subjects or elements about which we care interested in making decisions (inferences). A data frame is a list of the members of the population.</p> <p>A sample is a subset of the population which is used to gain insight about the population. Samples are used to represent the larger population of interest.</p> <p>A parameter is a numerical descriptive measure of a population. It is usually unknown as the population is too large to be collected. A population can have many parameters, and the parameters can be inferred (or estimated) from a sample (mostly random sample).</p>
Distributions	<p><u>Two types of distribution functions:</u></p> <ol style="list-style-type: none"> 1. Probability density function (for the continuous random variable) or probability mass function (for the discrete random variables). We denote them by pdf or pmf. They describe the probability (or likelihood) of the random variable at a particular value. 2. Cumulative probability function (cdf) which is the integral of the probability density function (pdf) or the cumulative sum of the probability mass function (pmf) on a given interval. You have computed p-values for testing the hypothesis, and the p-values are the application of cumulative probability function. <p><u>Bernoulli distribution:</u></p> <p>Bernoulli distribution is the discrete probability distribution of a random variable X which takes the value 1 with probability p and the value 0 with probability 1 – p. For example, toss a coin once. Let X = 1 denote the head and X = 0 the tail. The probability mass function (pmf) of Bernoulli distribution is:</p> $f_X(x) = p^x(1 - p)^{1-x}, x \in \{0, 1\} \text{ and } p = 0.5$ for tossing a coin once. <p>We can simply denote $X \sim \text{Bernoulli}(p)$</p> $E(X) = \sum x \cdot f_X(x)$ for discrete random variables $E(X) = \int x f_X(x) dx$ for continuous random variables <p><u>Binomial distribution:</u></p> <p>Binomial distribution describes the probability of getting exactly X = k successes (with probability p for each trial) in n independent Bernoulli trials. For example, we toss the coin 10 times and get 4 heads. The probability mass function (pmf) of Binomial distribution follows:</p> $f_X(k) = \binom{n}{k} p^k (1 - p)^{n-k}, k \in \{0, 1, 2, \dots, n\} \quad \# \binom{n}{k} = \frac{n!}{k!(n-k)!}$ <p>The cumulative distribution function can be expressed as:</p>

$$F(k; n, p) = \Pr(X \leq k) = \sum_{i=0}^k \binom{n}{i} p^i (1-p)^{n-i}, k \in \{0, 1, 2, \dots, n\}$$

We can denote $X \sim \text{Binomial}(n, p)$

Normal distribution:

Normal (or Gaussian) distributions are a type of continuous probability distribution for a real-valued random variable. Normal distributions play a central role in statistical inference!

The pdf for a normal random variable X is:

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), x \in (-\infty, \infty)$$

The cdf for a normal random variable X is:

$$F_X(x) = P(X \leq x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(z-\mu)^2}{2\sigma^2}\right) dz$$

We denote $X \sim N(\mu, \sigma^2)$

Z-score:

$$Z = \frac{X-\mu}{\sigma}$$

Central Limit Theorem:

Central Limit Theorem (CLT) Under some conditions, when the sample size n is large ($n \rightarrow \infty$), the sample statistic converges to a population parameter (This is a “dirty” definition of CLT).

t-distribution:

1. Degree of freedom (df) = the maximum number of logically independent values x_1, x_2, x_3 take any values, then $df = 3$. Now $x_1 + x_2 + x_3 = 2$, $df = 3 - 1 = 2$.
2. If $n \geq 30$, t-distribution approaches to $N(0, 1)$, so $N(0, 1)$ is used to approximate t-distribution for $n \geq 30$.
3. t-distribution are symmetric about zero, with longer tails than $N(0,1)$.
4. If a random value $X \sim t(df = k)$, then $X^2 \sim F(1, df = k)$.

Exact Binomial Test

Exact test, no need assumption check.

Step 1: formulate the hypothesis # k is the proportion

$H_0: p = k$ v.s. $H_1: p \neq k$ (two – tailed)

$H_1: p > k$ (right, one – sided, upper – tail)

$H_1: p < k$ (left, one – sided, lower – tail)

Step 2: compute the test statistics

$X^* = \text{number we observed}$

	<p>Step 3: compute the p-value or rejection region (RR) $H_1: p \neq k \Rightarrow p - \text{value} = 2 * P(X \geq X^*) = 2(1 - P(X \leq X^* - 1)) = 2(1 - \text{binomcdf}(n, k, X^* - 1))$ $H_1: p > k \Rightarrow p - \text{value} = P(X \geq X^*) = 1 - P(X \leq X^* - 1) = 1 - \text{binomcdf}(n, k, X^* - 1)$ $H_1: p < k \Rightarrow p - \text{value} = P(X \leq X^*) = \text{binomcdf}(n, k, X^*)$</p> <p>Step 4: Make conclusion. If p-value < α, we reject the null hypothesis.</p>
Two-sample Proportion Z-test	<p>Assumption:</p> <ol style="list-style-type: none"> 1) The number of success and failure must be greater than 5 for both samples. 2) Each sample was randomly taken from the population. <p>Step 1: formulate the hypothesis $H_0: p_1 - p_2 = 0$ v.s. $H_1: p_1 - p_2 \neq 0$ (two - sided) $H_1: p_1 - p_2 > 0$ (right, one - sided, upper - tail) $H_1: p_1 - p_2 < 0$ (left, one - sided, lower - tail)</p> <p>Step 2: Compute the test statistics Z^* based on CLT $Z^* = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\hat{p}(1 - \hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \sim N(0, 1)$, where $\hat{p}_1 = \frac{x_1}{n_1}$, $\hat{p}_2 = \frac{x_2}{n_2}$, $\hat{p} = \frac{x_1 + x_2}{n_1 + n_2}$</p> <p>Step 3: Calculate the p-value $H_1: p_1 - p_2 \neq 0 \Rightarrow p - \text{value} = 2 * \text{normalcdf}(Z^* , 10^{99}, 0, 1)$ $H_1: p_1 - p_2 > 0 \Rightarrow p - \text{value} = \text{normalcdf}(Z^*, 10^{99}, 0, 1)$ $H_1: p_1 - p_2 < 0 \Rightarrow p - \text{value} = \text{normalcdf}(Z^*, 10^{99}, 0, 1)$</p> <p>Step 4: Make conclusion. If p-value < α, we reject the null hypothesis.</p>
CI for two-sample Proportion Z-test	<p>$(1 - \alpha)100\%$ confidence interval for $p_1 - p_2$ is: $(\hat{p}_1 - \hat{p}_2) \pm z_{1 - \frac{\alpha}{2}} \sqrt{\hat{p}(1 - \hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$</p> <p>Interpretation: We are $(1 - \alpha)100\%$ confident that XXXX is between LOWER BOUND and UPPER BOUND.</p>
Two-sample t-test (Paired t-test, General)	<p>Assumption:</p> <ol style="list-style-type: none"> 1) The differences for the matched pairs, follow a normal probability distribution or the sample size (number of pairs) should be larger than or equal to 30. 2) The sample of pairs is a simple random sample from its population. <p>Step 1: formulate the hypothesis we define $\mu_1 - \mu_2 = d$ or $\mu_2 - \mu_1 = d$ $H_0: d = 0$ v.s. $H_1: d \neq 0$ (two - sided)</p>

	<p> $H_0: d \leq 0$ v.s. $H_1: d > 0$ (right, one – sided, upper – tail) $H_0: d \geq 0$ v.s. $H_1: d < 0$ (left, one – sided, lower – tail) </p> <p>Step 2: Compute the test statistics t^* based on CLT</p> $t^* = \frac{\bar{d}}{\hat{\sigma}_d/\sqrt{n}} \sim t(df = n - 1), \text{ where } \bar{d} = \frac{\sum d_i}{n}, \hat{\sigma}_d = \sqrt{\frac{\sum (d_i - \bar{d})^2}{n-1}}$ <p>Step 3: Calculate the p-value</p> <p> $H_1: d \neq 0 \Rightarrow p\text{-value} = 2 * tcdf(t^* , 10^{99}, n - 1)$ $H_1: d > 0 \Rightarrow p\text{-value} = tcdf(t^*, 10^{99}, n - 1)$ $H_1: d < 0 \Rightarrow p\text{-value} = tcdf(t^*, 10^{99}, n - 1)$ </p> <p>Step 4: Make conclusion. If p-value < α, we reject the null hypothesis.</p>
<p>CI for two-sample t-test (Paired t-test, General)</p>	<p>(1-α)100% confidence interval for the difference between (paired) two means is:</p> $\bar{d} \pm t(1 - \frac{\alpha}{2}, df = n - 1) \frac{\hat{\sigma}_d}{\sqrt{n}}$ <p>Interpretation: We are (1-α)100% confident that XXXX is between LOWER BOUND and UPPER BOUND.</p>
<p>Two-sample t-test ($\sigma_1^2 = \sigma_2^2$) ANOVA</p>	<p>Assumption:</p> <ol style="list-style-type: none"> 1) Two samples must be independent 2) Two samples are randomly generated from its population, respectively 3) Two populations have equal variance 4) Each sample size has more than 30 observations <p>Step 1: formulate the hypothesis # same with paired t-test</p> <p> $H_0: \mu_1 - \mu_2 = 0$ v.s. $H_1: \mu_1 - \mu_2 \neq 0$ (two – sided) $H_0: \mu_1 - \mu_2 \leq 0$ v.s. $H_1: \mu_1 - \mu_2 > 0$ (right, one – sided, upper – tail) $H_0: \mu_1 - \mu_2 \geq 0$ v.s. $H_1: \mu_1 - \mu_2 < 0$ (left, one – sided, lower – tail) </p> <p>Step 2: Compute the test statistics t^* based on CLT</p> $t^* = \frac{\bar{Y}_1 - \bar{Y}_2}{\hat{\sigma}_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t(df = n_1 + n_2 - 2)$ <p>where, $\hat{\sigma}_p = \sqrt{\frac{(n_1-1)\hat{\sigma}_1^2 + (n_2-1)\hat{\sigma}_2^2}{n_1+n_2-2}}$, $\hat{\sigma}_1^2 = \frac{1}{n_1-1} \sum_{i=1}^{n_1} (Y_{i1} - \bar{Y}_1)^2$, $\hat{\sigma}_2^2 = \frac{1}{n_2-1} \sum_{i=1}^{n_2} (Y_{i2} - \bar{Y}_2)^2$</p> <p>Step 3: Compute the p-value</p> <p> $H_1: \mu_1 - \mu_2 \neq 0 \Rightarrow p\text{-value} = 2 * tcdf(t^* , 10^{99}, n_1 + n_2 - 2)$ $H_1: \mu_1 - \mu_2 > 0 \Rightarrow p\text{-value} = tcdf(t^*, 10^{99}, n_1 + n_2 - 2)$ $H_1: \mu_1 - \mu_2 < 0 \Rightarrow p\text{-value} = tcdf(t^*, 10^{99}, n_1 + n_2 - 2)$ </p>

	<p>Step 4: Make conclusion. If p-value < α, we reject the null hypothesis.</p>
CI for two-sample t-test ($\sigma_1^2 = \sigma_2^2$) ANOVA	<p>(1-α)100% confidence interval for $\mu_1 - \mu_2$ is:</p> $(\bar{Y}_1 - \bar{Y}_2) \pm t \left(1 - \frac{\alpha}{2}, df = n_1 + n_2 - 2\right) \hat{\sigma}_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$ <p>Interpretation: We are (1-α)100% confident that XXXX is between LOWER BOUND and UPPER BOUND.</p>
Two-sample t-test ($\sigma_1^2 \neq \sigma_2^2$) NOT ANOVA	<p>Assumption:</p> <ol style="list-style-type: none"> 1) Two samples must be independent 2) Two samples are randomly generated from its population, respectively 3) Each sample size has more than 30 observations <p>Step 1: formulate the hypothesis # same with paired t-test</p> <p>$H_0: \mu_1 - \mu_2 = 0$ v.s. $H_1: \mu_1 - \mu_2 \neq 0$ (two - sided)</p> <p>$H_0: \mu_1 - \mu_2 \leq 0$ v.s. $H_1: \mu_1 - \mu_2 > 0$ (right, one - sided, upper - tail)</p> <p>$H_0: \mu_1 - \mu_2 \geq 0$ v.s. $H_1: \mu_1 - \mu_2 < 0$ (left, one - sided, lower - tail)</p> <p>Step 2: Compute the test statistics t^* based on CLT</p> $t^* = \frac{\bar{Y}_1 - \bar{Y}_2}{\sqrt{\frac{\hat{\sigma}_1^2}{n_1} + \frac{\hat{\sigma}_2^2}{n_2}}} \sim t(df = v)$ <p>where, $\hat{\sigma}_1^2 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (Y_{i1} - \bar{Y}_1)^2$, $\hat{\sigma}_2^2 = \frac{1}{n_2 - 1} \sum_{i=1}^{n_2} (Y_{i2} - \bar{Y}_2)^2$, $v = \frac{(n_1 - 1)(n_2 - 1)}{(1 - c)^2(n_1 - 1) + c^2(n_2 - 1)}$, $c = \frac{\frac{\hat{\sigma}_1^2}{n_1}}{\frac{\hat{\sigma}_1^2}{n_1} + \frac{\hat{\sigma}_2^2}{n_2}}$</p> <p>Step 3: Compute the p-value</p> <p>$H_1: \mu_1 - \mu_2 \neq 0 \Rightarrow p - value = 2 * tcdf(t^* , 10^{99}, v)$</p> <p>$H_1: \mu_1 - \mu_2 > 0 \Rightarrow p - value = tcdf(t^*, 10^{99}, v)$</p> <p>$H_1: \mu_1 - \mu_2 < 0 \Rightarrow p - value = tcdf(t^*, 10^{99}, v)$</p> <p>Step 4: Make conclusion. If p-value < α, we reject the null hypothesis.</p>
CI for two-sample t-test ($\sigma_1^2 \neq \sigma_2^2$) NOT ANOVA	<p>(1-α)100% confidence interval for $\mu_1 - \mu_2$ is:</p> $(\bar{Y}_1 - \bar{Y}_2) \pm t \left(1 - \frac{\alpha}{2}, df = v\right) \sqrt{\frac{\hat{\sigma}_1^2}{n_1} + \frac{\hat{\sigma}_2^2}{n_2}}$ <p>Interpretation: We are (1-α)100% confident that XXXX is between LOWER BOUND and UPPER BOUND.</p>
One-way ANOVA	<p>Assumption:</p> <ol style="list-style-type: none"> 1) Two samples must be independent 2) Two samples are randomly generated from its population, respectively 3) Two populations have equal variance

4) Each sample size has more than 30 observations

The setup of one-way ANOVA

The single-factor mean model is

$$Y_{ij} = \mu_1 + \epsilon_{ij}, \quad i = 1, 2, \dots, a, \quad j = 1, 2, \dots, n_i$$

Singe-factor effect model is

$$Y_{ij} = \mu + \tau_i + \epsilon_{ij}, \quad i = 1, 2, \dots, a, \quad j = 1, 2, \dots, n_i$$

Facts:

- 1) $i = 1, 2, \dots, a$ denotes the level of the factor/treatment. There are a levels.
- 2) $j = 1, 2, \dots, n_i$ denotes the subject in each factor/treatment level. There are n_1, n_2, \dots, n_a subjects in each level of the factor/treatment.
- 3) If $n_1 = n_2 = \dots = n_a$ (every level of the factor/treatment has the same number of subjects), it is called a balanced design. If not, it is an unbalanced design.
- 4) The mean model and effect model are equivalent, because $\mu_i = \mu + \tau_i$ for $i = 1, 2, \dots, a$.
- 5) ϵ_{ij} 's are model errors. They are unknown and estimated via residuals. The "variance" in Analysis of Variance means "the variance of the residuals".
- 6) Here we assume ϵ_{ij} 's are independent and identically distributed (IID) and $\epsilon_{ij} \sim N(0, \sigma^2)$ for all i and j .
- 7) The ANOVA model is called Fixed Effect, because the treatment effects are treated as unknown constants (fixed values) instead of random values. If they are treated as random values, it is called a random effect model.

Notations:

$$\hat{\mu} = \bar{Y} \text{ (Grand mean), } E(\hat{\mu}) = \mu$$

$$\hat{\mu}_i = \bar{Y}_i \text{ (Mean in each level), } E(\hat{\mu}_i) = \mu_i$$

$$\hat{\tau}_i = \bar{Y}_i - \bar{Y}_., \text{ } E(\hat{\tau}_i) = \tau_i$$

Calculation:

Total sample size: $N = n_1 + n_2 + \dots + n_a$

Overall (Grand) mean (estimate for μ): $\bar{Y}_. = \frac{1}{n} \sum_{i=1}^a \sum_{j=1}^{n_i} Y_{ij}$

Treatment means (estimates for μ_i 's): $\bar{Y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}, \quad i = 1, 2, \dots, a$

Treatment effects (estimates for τ_i 's): $\bar{Y}_i - \bar{Y}_., \quad i = 1, 2, \dots, a$

$$SS_T = \sum_{i=1}^a \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_.)^2, \quad df(SS_T) = N - 1$$

$$SS_{trt} = \sum_{i=1}^a n_i (\bar{Y}_i - \bar{Y}_.)^2, \quad df(SS_{trt}) = df(MS_{trt}) = a - 1, \quad MS_{trt} = \frac{SS_{trt}}{a-1}$$

$$SS_E = \sum_{i=1}^a \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2, \quad df(SS_E) = df(MS_E) = N - a, \quad MS_E = \frac{SS_E}{N-a}$$

$$SS_T = SS_{trt} + SS_E$$

$$E(MS_E) = \sigma^2$$

$$E(MS_{trt}) = \sigma^2 + \frac{\sum_{i=1}^a n_i \tau_i^2}{a-1}$$

Conducting F-test

Step 1: formulate the hypothesis

$$H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_a, \quad H_1: \mu_i \neq \mu_j \text{ for at least one pair of } (i, j)$$

or

$$H_0: \tau_1 = \tau_2 = \tau_3 = \dots = \tau_a, \quad H_1: \tau_i \neq 0 \text{ for at least one } i$$

Step 2: Compute the test statistics F^*

$$F^* = \frac{MS_{trt}}{MS_E} \sim F(df_1 = a - 1, df_2 = N - a)$$

Here we use the upper tail, so the test statistics here is always greater than 1

df_1 is the degree of freedom of numerator (MS_{trt}) and df_2 is the degree of freedom of denominator (MS_E)

It is okay to switch the numerator and the denominator, but the degrees of freedom and p-value of the F distribution should also be switched.

Step 3: Compute the p-value

$$p\text{-value} = Fcdf(F^*, 10^{99}, df_1 = a - 1, df_2 = N - a)$$

Step 4: Make conclusion. If p-value < α , we reject the null hypothesis.

ANOVA table

Source of Variation	Sum of Squares	df	Mean Square	F-score
Between Treatments	SS_{trt}	$a - 1$	$MS_{trt} = \frac{SS_{trt}}{a - 1}$	$F^* = \frac{MS_{trt}}{MS_E}$
Error(Within Treatments)	SS_E	$N - a$	$MS_E = \frac{SS_E}{N - a}$	
Total	SS_T	$N - 1$		

Model Assumption Check

1. The sample (the design) is a random sample.

2. The errors are normally distributed

Method 1: QQ plot (quantile-quantile plot).

Method 2: Shapiro-Wilk normality test.

H_0 : The data are normally distributed, H_1 : The data are NOT normally distributed
 If the p-value $< \alpha$, we reject the null hypothesis.

3. The errors are independently distributed

Method 1: Plot of residuals in time series (if the order of the runs is known)

Method 2: Plot of residuals $\hat{\epsilon}_{ij}$'s versus fitted values \bar{Y}_i 's.

4. The errors have equal variances

Method 1: Grouped boxplot.

If the widths of the boxes are similar, then we conclude they have equal variance. (NOT accurate)

Method 2: Tests for equal variance

H_0 : All groups have equal variances ($\sigma_1^2 = \sigma_2^2 = \dots = \sigma_a^2$), H_1 : All groups have UNEQUAL variances

Levene's test: Assess the equality of variances for a variable calculated for two or more samples (groups/levels).

1) $d_{ij} = |Y_{ij} - \bar{Y}_i|$, for $i = 1, 2, \dots, a$; $j = 1, 2, \dots, n_i$

2) Run on-way ANOVA test on d_{ij}

Brown-Forsythe test: Use median instead of mean in Levene's test.

$d_{ij} = |Y_{ij} - \tilde{Y}_i|$, for $i = 1, 2, \dots, a$; $j = 1, 2, \dots, n_i$

Bartlett's test: Derived from the likelihood ratio test under normality assumption. (NOT required for this class)

Brown-Forsythe test should be preferred, it is the "median" version of Levene's test. If the data are skewed, Brown-Forsythe is more robust; if the data are symmetric (not skewed), Brown-Forsythe test and Levene's test have similar performance.

Comparisons among treatment level means

$H_0: \mu_i = \mu_j$, $H_1: \mu_i \neq \mu_j$ or $H_0: \mu_{i_1} + \mu_{i_2} = \mu_{j_1} + \mu_{j_2}$, $H_1: \mu_{i_1} + \mu_{i_2} \neq \mu_{j_1} + \mu_{j_2}$

We can use contrast to perform this above

$$\Gamma = \sum_{i=1}^a c_i \mu_i$$

$$H_0: \sum_{i=1}^a c_i \mu_i = 0, \quad H_1: \sum_{i=1}^a c_i \mu_i \neq 0$$

$$\text{Point estimate: } C = \sum_{i=1}^a c_i \bar{Y}_i.$$

$$\text{Sample variance: } \text{Var}(C) = MS_E \sum_{i=1}^a \frac{c_i^2}{n_i}$$

$$\text{Standard error: } s(C) = \sqrt{\text{Var}(C)} = \sqrt{MS_E \sum_{i=1}^a \frac{c_i^2}{n_i}}$$

$$\text{Test statistic: } F^* = \frac{C^2}{\text{Var}(C)} = \frac{(\sum_{i=1}^a c_i \bar{Y}_i)^2}{MS_E \sum_{i=1}^a \frac{c_i^2}{n_i}} \sim F(df_1 = 1, df_2 = N - a)$$

We reject the null if p-value $< \alpha$.

	<p>(1-α)100% confidence interval on the contrast $\Gamma = \sum_{i=1}^a c_i \mu_i$ is: $C \pm t \left(1 - \frac{\alpha}{2}, df = N - a \right) s(C)$</p> <p>Scheffé Method for Comparing all contrasts (simultaneous method)</p> $s(C) = \sqrt{Var(C)} = \sqrt{MS_E \sum_{i=1}^a \frac{c_i^2}{n_i}}$ <p>We reject the null hypothesis $H_0: \sum_{i=1}^a c_i \mu_i = 0$ if $C > s(C) \sqrt{(a-1)F(1-\alpha, df_1 = a-1, df_2 = N-a)}$ Simultaneous confidence intervals: $C \pm s(C) \sqrt{(a-1)F(1-\alpha, df_1 = a-1, df_2 = N-a)}$</p> <p>Comparing pairs of treatment means $H_0: \mu_i - \mu_j = 0, H_1: \mu_i - \mu_j \neq 0$ for all $i \neq j$</p> <p>Tukey HSD: We reject null hypothesis if: $\bar{Y}_i - \bar{Y}_j > \frac{q_{\alpha}(a, df=N-a)}{\sqrt{2}} \sqrt{MS_E \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}$ Tukey's (1-α)100% confidence interval for $\mu_i - \mu_j$ is $(\bar{Y}_i - \bar{Y}_j) \pm \frac{q_{\alpha}(a, df=N-a)}{\sqrt{2}} \sqrt{MS_E \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}$</p> <p>Fisher Least Significant Difference (LSD): We reject null hypothesis if: $\bar{Y}_i - \bar{Y}_j > t \left(1 - \frac{\alpha}{2}, df = N - a \right) \sqrt{MS_E \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}$ Fisher's (1-α)100% confidence interval for $\mu_i - \mu_j$ is $(\bar{Y}_i - \bar{Y}_j) \pm t \left(1 - \frac{\alpha}{2}, df = N - a \right) \sqrt{MS_E \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}$ # same as CI on the difference in any two treatment means $\mu_i - \mu_j$</p>
CI for one-way ANOVA (pairwise)	<p>(1-α)100% confidence interval for the i-th treatment mean μ_i is: $\bar{Y}_i \pm t \left(1 - \frac{\alpha}{2}, df = N - a \right) \sqrt{\frac{MS_E}{n_i}}$</p> <p>(1-$\alpha$)100% confidence interval on the difference in any two treatment means $\mu_i - \mu_j$ is: $(\bar{Y}_i - \bar{Y}_j) \pm t \left(1 - \frac{\alpha}{2}, df = N - a \right) \sqrt{MS_E \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}$</p>