# Assumptions of the ANOVA F-test:

• Again, most assumptions involve the $\varepsilon_{ij}$'s (the error terms).

(1) The model is correctly specified.
(2) The $\varepsilon_{ij}$'s are normally distributed.
(3) The $\varepsilon_{ij}$'s have mean zero and a common variance, $\sigma^2$.
(4) The $\varepsilon_{ij}$'s are independent across observations.

• With multiple populations, detection of violations of these assumptions requires examining the residuals rather than the $Y$-values themselves.

• An estimate of $\varepsilon_{ij}$ is: $Y_{ij} - \hat{\mu}_{i}$

$$= Y_{ij} - \bar{Y}_{i\cdot}$$

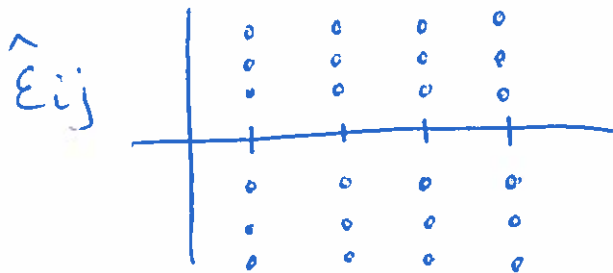• Hence the residual for data value $Y_{ij}$ is: $Y_{ij} - \bar{Y}_{i\cdot}$

• We can check for non-normality or outliers using residual plots (and normal Q-Q plots) from the computer.

• Checking the equal-variance assumption may be done with a formal test:
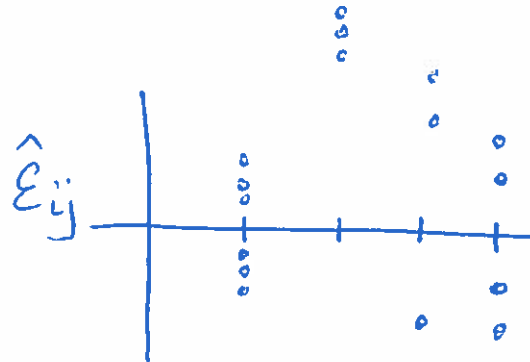
$H_0: \sigma_1^2 = \sigma_2^2 = \ldots = \sigma_t^2$
$H_a:$ at least two variances are not equal

- **The Levene test is a formal test for unequal variances that is robust to the normality assumption.**

- **It performs the ANOVA F-test on the absolute residuals from the sample data.**

**Example pictures:**

$\hat{\varepsilon}_{ij}$

$\hat{\varepsilon}_{ij}$

Levene test : Won't reject $H_0$

Levene test **will** probably reject $H_0$.

- For Rice data: Levene test has p-value $= 0.4654 > .05$. Conclude the equal-variances assumption is reasonable.
- Normal Q-Q plot shows the normality assumption is a bit questionable.

## Remedies to Stabilize Variances

- **If the <u>variances appear unequal</u> across populations, using transformed values of the response may remedy this. (Such transformations can also help with violations of the <u>normality assumption</u>.)**

- **The drawback is that interpretations of results may be less convenient.**

Suggested transformations:

- If the standard deviations of the groups increase proportionally with the group means, try: $Y_{ij}^* = \log(Y_{ij})$
- If the variances of the groups increase proportionally with the group means, try: $Y_{ij}^* = \sqrt{Y_{ij}}$
- If the responses are proportions (or percentages), try: $Y_{ij}^* = \arcsin(\sqrt{Y_{ij}})$

- If none of these work, may need to use a nonparametric procedure (e.g., Kruskal-Wallis test).

## Making Specific Comparisons Among Means

- If our F-test rejects $H_0$ and finds there are significant differences among the population means, we typically want more specific answers:

(1) Is the mean response at a specified level superior to (or different from) the mean response at other levels?

(2) Is there some natural grouping or separation among the factor level mean responses?

- Question (1) involves a "pre-planned" comparison and is tested using a contrast.

- Question (2) is a "post-hoc" comparison and is tested via a "Post-Hoc Multiple Comparisons" procedure.

# Contrasts

- **A contrast is a linear combination of the population means whose coefficients add up to zero.**

**Example ($t = 4$):** $4\mu_1 + 7\mu_2 - 13\mu_3 + 2\mu_4$

- **Often a contrast is used to test some meaningful question about the mean responses.**

**Example (Rice data):  Is the mean of variety 4 different from the mean of the other three varieties?**

**We are testing:** $H_0: \dfrac{\mu_1 + \mu_2 + \mu_3}{3} = \mu_4$

vs. $H_a: \dfrac{1}{3}\mu_1 + \dfrac{1}{3}\mu_2 + \dfrac{1}{3}\mu_3 \neq \mu_4$

**What is the appropriate contrast?**

$L = \dfrac{1}{3}\mu_1 + \dfrac{1}{3}\mu_2 + \dfrac{1}{3}\mu_3 - \mu_4$  (coefficients add to zero)

**Now we test:** $H_0: L = 0$
$H_a: L \neq 0$

**We can estimate $L$ by:**

$\hat{L} = \dfrac{1}{3}\bar{Y}_{1\cdot} + \dfrac{1}{3}\bar{Y}_{2\cdot} + \dfrac{1}{3}\bar{Y}_{3\cdot} - \bar{Y}_{4\cdot}$

**Under $H_0$, and with balanced data, the variance of a contrast**  $\hat{L} = a_1\bar{Y}_{1\cdot} + \cdots + a_t\bar{Y}_{t\cdot}$

**is:** $\mathrm{var}\left(\hat{L}\right) = \left(a_1^2 + \cdots + a_t^2\right)\dfrac{\sigma^2}{n}$

- **Also, when the data come from normal populations, $\hat{L}$ is normally distributed.**

- **Replacing $\sigma^2$ by its estimate MSW:**

$$t^* = \frac{\hat{L}}{\sqrt{\widehat{var}(\hat{L})}}$$

has a t-distribution under $H_0$ with $df = t(n-1)$ (assuming $n_1 = \cdots = n_t = n$)

**For balanced data:**

$$t^* = \frac{\sum_i a_i \bar{Y}_{i\bullet}}{\sqrt{\frac{MSW}{n}\sum_i a_i^2}}$$

- **To test $H_0$: $L = 0$, we compare $t^*$ to the appropriate critical value in the t-distribution with $t(n-1)$ d.f.**

- **Our software will perform these tests even if the data are unbalanced.**

$$L = \tfrac{1}{3}\mu_1 + \tfrac{1}{3}\mu_2 + \tfrac{1}{3}\mu_3 - \mu_4$$

$\alpha = .05$  **Example:** Test $H_0 : L = 0$  vs.  $H_a : L \neq 0$

$$t^* = \frac{-166.0833}{37.221} = -4.46$$

Compare $|t^*|$ to $t_{.025}(12\ d.f.) = 2.179$

$|t^*| = 4.46 > 2.179$, and also P-value $= .0008 < .05$, so we reject $H_0$. Conclude mean yield for variety 4 differs from mean yield of other varieties.

- **Note: When testing multiple contrasts, the specified $\alpha$ (= P{Type I error} ) applies to each test individually, not to the <u>series</u> of tests collectively.**

Example 2:  $L = \mu_1 - \mu_2$    $H_0 : L = 0$
                                                                $H_a : L \neq 0$

$\Rightarrow$ P-value $= .2409 \rightarrow$ fail to reject $H_0$

# Post Hoc Multiple Comparisons

- **When we specify a significance level $\alpha$, we want to limit P{Type I error}.**

- **What if we are doing many simultaneous tests?**

- **Example: We have $\mu_1$, $\mu_2$, ..., $\mu_t$. We want to compare <u>all pairs</u> of population means.**

- **<u>Comparisonwise error rate</u>: The probability of a Type I error on <u>each comparison</u>.**

- **<u>Experimentwise error rate</u>: The probability that the simultaneous testing results in <u>at least one</u> Type I error.**

- **We only do post hoc multiple comparisons if the overall F-test indicates a difference among population means.**
- **If so, our question is: Exactly <u>which</u> means are different?**

*actually a series of null hypotheses.*

- **We test:** $H_0: \mu_i = \mu_j$ for <u>all</u> $i \neq j$

- **The <u>Fisher LSD procedure</u> performs a t-test for each pair of means (using a common estimate of $\sigma^2$, MSW).**

- **The Fisher LSD procedure declares $\mu_i$ and $\mu_j$ significantly different if:**

$$\left| \bar{Y}_{i\cdot} - \bar{Y}_{j\cdot} \right| > t_{\alpha/2} \sqrt{\frac{2\,MSW}{n}}$$

$df = $ "within-groups d.f."

$\uparrow$ assuming balanced data

- **Problem:** Fisher LSD only controls the <u>comparisonwise</u> error rate.
- The <u>experimentwise</u> error rate may be <u>much larger</u> than our specified $\alpha$.

- <u>Tukey's Procedure</u> controls the <u>experimentwise</u> error rate to be only equal to $\alpha$.

- Tukey procedure declares $\mu_i$ and $\mu_j$ significantly different if:

$$\left| \bar{Y}_{i\cdot} - \bar{Y}_{j\cdot} \right| > q_\alpha(t, df) \sqrt{\frac{MSW}{n}} \leftarrow \text{balanced data}$$

- $q_\alpha(t, df)$ is a critical value based on the studentized range of sample means:

$$q = \frac{\left( \bar{Y}_{max} - \bar{Y}_{min} \right)}{\sqrt{MSW/n}}$$

- Tukey critical values are listed in Table A.7.

- Note: $q_\alpha(t, df)$ is larger than $\sqrt{2} \left( t_{\alpha/2} \right)$

$\rightarrow$ Tukey procedure will declare a significant difference between two means __less__ often than Fisher LSD.

$\rightarrow$ Tukey procedure will have __lower__ experimentwise error rate, but Tukey will have __less__ power than Fisher LSD.

$\rightarrow$ Tukey procedure is a __more__ conservative test than Fisher LSD.

## Some Specialized Multiple Comparison Procedures

- **Duncan multiple-range test**: An adjustment to Tukey's procedure that reduces its conservatism.
- **Dunnett's test**: For comparing several treatments to a "control".
- **Scheffe's procedure**: For testing "all possible contrasts" rather than just all possible pairs of means.

**Notes**: - **When appropriate**, preplanned comparisons are considered superior to post hoc comparisons (more power).

- **Tukey's procedure can produce simultaneous CIs for all pairwise differences in means.** Produces CIs for $M_i - M_j$ for <u>all</u> $i \neq j$
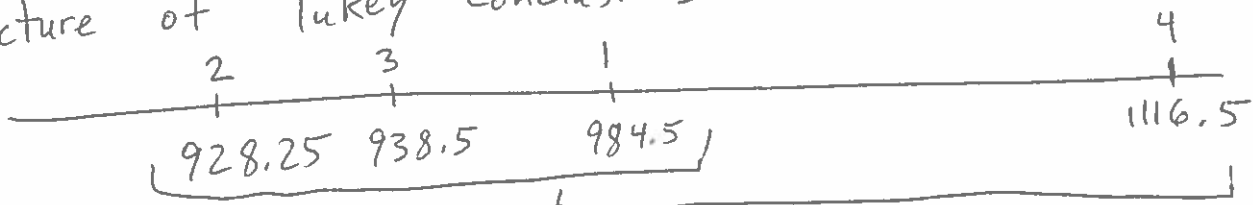
**Example**: Rice data:

Fisher LSD (using $\alpha = .05$) declares:

$M_1$ and $M_4$  are  significantly different

$M_2$ and $M_4$  "  "  "

$M_3$ and $M_4$  "  "  "

---

Tukey (using $\alpha = .05$) declares:

$M_2$ and $M_4$  are  signif. different

$M_3$ and $M_4$  "  "  "

Picture of Tukey conclusions:

```
        2     3        1                    4
    ----+-----+--------+--------------------+----
      928.25 938.5   984.5               1116.5
```

# Random Effects Model

$$H_0: \mu_1 = \cdots = \mu_t$$

$$\Updownarrow$$

$$H_0: \tau_1 = \cdots = \tau_t = 0$$

- **Recall our ANOVA model:**

$$Y_{ij} = \mu_i + \varepsilon_{ij}, \quad i = 1, \ldots, t, \quad j = 1, \ldots, n_i$$

$$Y_{ij} = \mu + \tau_i + \varepsilon_{ij}$$

- **If the *t* levels of our factor are the only levels of interest to us, then $\tau_1, \tau_2, \ldots, \tau_t$ are called <u>fixed effects</u>.**

- **If the *t* levels represent a random selection from a <u>large population</u> of levels, then $\tau_1, \tau_2, \ldots, \tau_t$ are called <u>random effects</u>.**

<u>Example</u>:  **From a population of teachers, we randomly select 6 teachers and observe the standardized test scores for their students.  Is there <u>significant variation</u> in student test score <u>among the population</u> of teachers?**

- **If $\tau_1, \tau_2, \ldots, \tau_t$ are random variables, the F-test no longer tests:**

$$H_0: \tau_1 = \tau_2 = \cdots = \tau_t = 0$$

**Instead, we test:**

$$H_0: \sigma_\tau^2 = 0$$

$$\text{vs.} \quad H_a: \sigma_\tau^2 > 0$$

<u>Question of interest</u>:  **Is there significant variation among the different levels in the population?**

↑ effects for the

- **For the one-way ANOVA, the test statistic is exactly the same, F\* = MSB / MSW, for the random effects model as for the fixed effects model.**