**Note 5** (Some notes about the $F$-test in ANOVA)

- Connection between the 2-sample (unpaired) $t$-test and the $F$-test in ANOVA ($I$ samples, $I \geq 2$): the 2-sample $t$-test is a special case of the $F$-test where only two groups are being compared ($I = 2$) because
  - test statistic (**recall**: LN, CH11, p.13, $t$-test for $n$ $X_i$'s and $m$ $Y_j$'s)

$J_1 \to \overline{Y_{1j}}$   $J_2 \to \overline{Y_{2j}}$

$*$ $SS_B/(I-1) = SS_B$      $\overline{Y_{..}}$   $\overline{Y_{2\cdot}}$

$\boxed{\overline{Y_{1\cdot}}}$

$\underset{J_1\to}{=} n\left[\overline{X} - \left(\frac{n}{m+n}\overline{X} + \frac{m}{m+n}\overline{Y}\right)\right]^2 + \underset{J_2\to}{m}\left[\overline{Y} - \left(\frac{n}{m+n}\overline{X} + \frac{m}{m+n}\overline{Y}\right)\right]^2$

$= n\frac{m^2}{(m+n)^2}(\overline{X}-\overline{Y})^2 + m\frac{n^2}{(m+n)^2}(\overline{X}-\overline{Y})^2 = \frac{mn}{m+n}(\overline{X}-\overline{Y})^2$

$*$ $SS_W/(N-I) = SS_W/[(m-1)+(n-1)] = s_p^2$ (LN, CH11, P.7, Definition 1)

$\boxed{\begin{array}{l}\text{2-sided} \\ |T| > c \\ \text{becomes} \\ \text{1-sided} \\ F > c'\end{array}}$

$\circledast$ $F = \dfrac{SS_B/(I-1)}{SS_W/(N-I)} = \dfrac{(\overline{X}-\overline{Y})^2}{s_p^2\left(\frac{1}{n}+\frac{1}{m}\right)} = T^2$

$\boxed{\text{indep.}}$   $\boxed{\chi_1^2/1}$

$t_{\underline{d}} = \dfrac{N(0,1)}{\sqrt{\chi_{\underline{d}}^2/d}}, \quad t_{\underline{d}}^2 = \dfrac{[N(0,1)]^2}{\chi_{\underline{d}}^2/d}$

$\boxed{\begin{array}{l}d = N-I \\ = m+n-2\end{array}}$   $\ominus$ null distribution: if $\underline{Z} \sim t_{\underline{d}}$, then $\underline{Z}^2 \sim F_{1,\underline{d}}$

- Under the model ($\circledast$) in LNp.5, the $F$-test in ANOVA is equivalent to the likelihood ratio test (exercise, the proof is similar to what presented in LN, CH11, p.15-18, for the case of two independent samples).

$\boxed{\begin{array}{l}\bullet \text{ Under } \Omega = H_0 \cup H_A, \text{ MLEs:} \\ \widehat{\mu}_{i,\Omega} = \overline{Y_{i\cdot}}, i = 1,2,\cdots,I, \\ \widehat{\sigma}_\Omega^2 = SS_W/N \\ \bullet \text{ Under } \omega = H_0, \text{ MLEs:} \\ \widehat{\mu}_{1,\omega} = \cdots = \widehat{\mu}_{I,\omega} = \overline{Y_{..}} \\ \widehat{\sigma}_\omega^2 = SS_{TOT}/N\end{array}}$

---

- **A nonparametric method --- the Kruskal-Wallis test** ← Why need it? Check Note 6, LN, CH11, p.25-26.

  - Consider the model ($\otimes$) in LNp.1, and further assume that
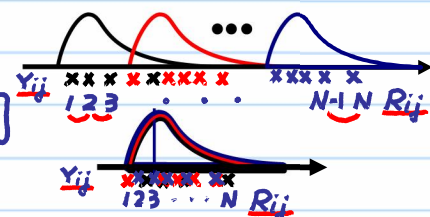    - $\Omega$ is the collection of *all* continuous distributions $\Rightarrow \dim(\Omega) = \infty$
    - $F_1, \ldots, F_I$ have the same shape, i.e.,

    $\boxed{\begin{array}{l}\text{Why need} \\ \text{this assumption?} \\ \text{Check Q.7} \\ \text{(LN, CH11, p35)}\end{array}}$   $F_1 = F(x-\Delta_1), \ldots, F_I = F(x-\Delta_I),$

    where $\underline{F} \in \Omega$.   $\boxed{F_1, \cdots, F_I \text{ have same variance}}$

  - Test the null and alternative hypotheses:

    $H_0 : \underline{\Delta_1} = \cdots = \underline{\Delta_I} = \underline{0}$   vs.   $H_A$ : at least one of $\Delta_i$'s is not $0$

$\bigstar \odot$ Under $H_0$, all $Y_{ij}$'s $\sim$ i.i.d. $F$. → Under $H_0$, the distribution of ranks is irrelevant to $F$.

- **Recall**. The sample size of the $i$th sample is $J_i$, $i = 1, \ldots, I$, and $N = J_1 + \cdots + J_I$ is the number of all observations.

$\boxed{\begin{array}{c}Y_{ij} \\ \downarrow \\ R_{ij}\end{array}}$ cf.

**Theorem 9** (Kruskal-Wallis test)   use ranks, rather than raw data, to do analysis

$\odot$ Let $R_{ij}$'s be the *ranks* of $Y_{ij}$'s in the *combined (pooled)* sample.

- Define   $\boxed{\text{weighted average of } \overline{R}_{1\cdot}, \cdots, \overline{R}_{I\cdot}}$

$\boxed{\overline{Y_{i\cdot}}}$ cf. $\ominus$ $\overline{R}_{i\cdot} = \dfrac{1}{J_i}\sum_{j=1}^{J_i} R_{ij}$: average rank in the $i$th group

$\boxed{\overline{Y_{..}}}$ cf. $\ominus$ $\overline{R}_{..} = \dfrac{1}{N}\sum_{i=1}^{I}\sum_{j=1}^{J_i} R_{ij} = \dfrac{J_1\overline{R}_{1\cdot} + \cdots + J_I\overline{R}_{I\cdot}}{J_1 + \cdots + J_I} = \dfrac{1}{N}\dfrac{N(N+1)}{2} = \dfrac{N+1}{2}$

$\underset{1+2+\cdots+N}{}$   $\boxed{\begin{array}{l}= 1 + \cdots + N \\ \text{a constant,} \\ \text{not r.v.}\end{array}}$

- Define

$$\boxed{Var(Z) = E(Z^2) - [E(Z)]^2, \; Z = \overline{R}_{i\cdot} \text{ with prob. } J_i/N}$$

$$\underline{SS_B} = \sum_{i=1}^{I} J_i (\overline{R}_{i\cdot} - \overline{R}_{\cdot\cdot})^2 = \left(\sum_{i=1}^{I} J_i \overline{R}_{i\cdot}^2\right) - N\overline{R}_{\cdot\cdot}^2 = \left(\sum_{i=1}^{I} J_i \overline{R}_{i\cdot}^2\right) - \frac{N(N+1)^2}{4}$$

**Between**

which measures *dispersion* of $\overline{R}_{i\cdot}$'s. — **large.** if $\Delta i$'s **very different**
                                      **Small.** if $\Delta i$'s **about the same**

- Test statistic $\underline{K}$

$$F \propto \frac{SS_B}{SS_W} \xleftarrow{\text{cf.}} \underline{K} = \frac{12}{N(N+1)} \underline{SS_B} = \frac{12}{N(N+1)} \left(\sum_{i=1}^{I} J_i \overline{R}_{i\cdot}^2\right) - 3(N+1)$$

(**Note.** $SS_B$ can be found by running $R_{ij}$'s through an ANOVA program)

  – **Q**: Why is $SS_B$ divided by $\frac{N(N+1)}{12}$ in $\underline{K}$? Define — $= 1^2 + 2^2 + \cdots + N^2 = \frac{N(N+1)(2N+1)}{6}$

$\boxed{Var(Z) = E(Z^2) - [E(Z)]^2,}$                                         $\boxed{\text{a constant, not } r.v.}$

$Z = R_{ij}$ with

$\text{prob. } 1/N$    $\underline{SS_{TOT}} = \sum_{i=1}^{I}\sum_{j=1}^{J_i}(R_{ij} - \overline{R}_{\cdot\cdot})^2 = \left(\sum_{i=1}^{I}\sum_{j=1}^{J_i} R_{ij}^2\right) - N\overline{R}_{\cdot\cdot}^2 = \frac{(N-1)\,N(N+1)}{12}$

Then,

$\sigma^2 \xleftarrow{\text{e }(H_0 \cup H_A)} \dfrac{SS_W}{N-I}$

       cf. $\updownarrow$     $\underline{K} = \dfrac{SS_B}{SS_{TOT}/(N-1)} = \dfrac{SS_B}{(SS_B + SS_W)/(N-1)} = \dfrac{1}{\left(1 + \frac{SS_W}{SS_B}\right)/(N-1)}$

$\sigma_R^2 \xleftarrow{\text{e (under } H_0)}$

**a constant**                    **Note.** $SS_{TOT} = SS_B + SS_W$ (LNp.7) still holds for $R_{ij}$'s.

$\Rightarrow$ $\boxed{\begin{array}{l} F \propto \frac{SS_B}{SS_W} \uparrow \\ \Leftrightarrow \frac{SS_B}{SS_{TOT}} \uparrow \\ \Leftrightarrow K \uparrow \end{array}}$

  – **Q**: Why is no $SS_W$ in $\underline{K}$?   $\boxed{\begin{array}{l}\text{a constant} \\ - SS_B\end{array}}$

  – Data with large values of $K$ are more extreme,
     i.e., provide stronger evidence against $H_0$.

$\boxed{\textbf{Q}: \text{Why ANOVA use} \; \boxed{\text{indep}} \frac{SS_B}{SS_W}, \text{not} \frac{SS_B}{SS_{TOT}} \boxed{\begin{array}{l}\text{not} \\ \text{indep}\end{array}}}$