overall F·test in LM → **F-Test**   noncentral parameter $\frac{\sum_i n_i \tau_i / N}{}$

$\frac{1}{k-1}\left[\|P_{\underline{n\Omega w}}\underline{\mu}\|^2 + (k-1)\sigma^2\right]$   By (P4) in LN p. 2-36 → $\sigma^2 + \frac{\sum_i n_i(\tau_i - \bar{\tau})^2}{k-1}$

ANOVA Table

| Source | Degrees of Freedom ($df$) | Sum of Squares | $RSS_w - RSS_\Omega = \|P_{\underline{n\Omega w}}\underline{y}\|^2$ | Mean Squares | Expected MS |
|---|---|---|---|---|---|
| between → treatment | $k-1$ | $SSTr = \sum_{i=1}^k n_i(\bar{y}_{i\cdot} - \bar{y}_{\cdot\cdot})^2$ | | $MSTr = SSTr/df$ | $E_\Omega(MSTr)$ |
| within → residual | $N-k$  $\mu_i's$ | $SSE = \sum_{i=1}^k \sum_{j=1}^{n_i}(y_{ij}-\bar{y}_{i\cdot})^2$  $RSS_\Omega$ | | $MSE = SSE/df = \hat{\sigma}^2$   (F) | $E_\Omega(MSE)$ |
| total | $N-1$ | $\sum_{i=1}^k \sum_{j=1}^{n_i}(y_{ij}-\bar{y}_{\cdot\cdot})^2$  $RSS_w$ | | By (P5) in LN p. 2-36 | $\sigma^2$ |

$\sum n_i$   $\tau$

noncentral parameter

The $F$ statistic for the <u>null hypothesis</u> that there is <u>no difference</u> between the <u>treatments</u>, i.e.,

$\underline{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_1 \\ \mu_2 \\ \vdots \\ \mu_2 \\ \vdots \\ \mu_k \\ \vdots \\ \mu_k \end{bmatrix}$   $\|P_{\underline{n\Omega w}}(\underline{\mu})\|^2 = \sum n_i(\mu_i - \bar{\mu})^2$

non-central parameter $= 0$ under $H_0$

$\mu_1 = \cdots = \mu_K \ (\mu_i = \tau + \tau_i)$

$H_0: \tau_1 = \cdots = \tau_k, \ = \tau$

$\underline{\mu} = \begin{matrix} m \\ \Omega \end{matrix}$   $= \underline{\mu} - \frac{\underline{\mu}^T \underline{1}}{\|\underline{1}\|^2}\underline{1}$

It follows a <u>noncentral chi-square</u> dist. under $\Omega\backslash\omega$ (LM, LN p. 4-7) & (LN p. 2-35)

$\Omega: y_{ij} = \tau + \tau_i + \varepsilon_{ij} \Rightarrow df_\Omega = N-k$

$\omega: y_{ij} = (\tau + \tau) + \varepsilon_{ij} \Rightarrow df_\omega = N-1$

is

$\bar{b} = \frac{\sum n_i \mu_i}{N} = \bar{\mu}$

Intuition: graphs in LN p.3-2~3

$\underline{F} = \frac{\sum_{i=1}^k n_i(\bar{y}_{i\cdot} - \bar{y}_{\cdot\cdot})^2/(k-1)}{\sum_{i=1}^k \sum_{j=1}^{n_i}(y_{ij}-\bar{y}_{i\cdot})^2/(N-k)} = \frac{MSTr}{MSE}$,

$F = \frac{(RSS_w - RSS_\Omega)/df_w - df_\Omega}{RSS_\Omega/df_\Omega}$

$\sim \sigma^2 \chi^2_{df_\Omega}$ under $\Omega$ & $\omega$

which has an $F$ distribution with parameters $k-1$ and $N-k$. (under $\omega$)

---

# ANOVA for Pulp Experiment

qualitative factor →

| Source | Degrees of Freedom ($df$) | Sum of Squares | Mean Squares | F |
|---|---|---|---|---|
| operator | 3 | 1.34 | 0.447 | 4.20 |
| residual | 16 | 1.70 | 0.106 | |
| total | 19 | 3.04 | | |

$\hat{\sigma}^2$

→ 20 obs.

- $Prob(F_{3,16} > 4.20) = \underline{0.02} = $ p-value,

  thus declaring a <u>significant</u> operator-to-operator <u>difference</u> at level 0.02.

  <u>Note</u>. ANOVA answers the question: "whether <u>there exist</u> <u>difference</u> between treatment means?"

- <u>Further question</u>: among the $6 = \binom{4}{2}$ pairs of operators, what pairs show significant difference?

  <u>Answer</u>: Need to use <u>multiple comparisons</u>. → answer "how different?" question.

# Constraint on the Parameters $\tau_i$'s

Identifiability (LM, LNp 5-10~11)

$\because \hat{\tau}_i\text{'s}, \hat{\eta}$ have infinite many solutions

- The model in LNp.3-4 has $k$ distinct levels, but $k+1$ regression parameters

$\because \omega, \Omega$ not change under the over-parameterized model

$\Rightarrow$ over-parameterized $\Rightarrow X^T X$ singular $\Rightarrow$ unidentifiable

$\Rightarrow$ cannot estimate parameters (**Q**: but why can do overall $F$-test?)

- Some common constraint on $\tau_i$'s ($\Omega: y_{ij} = \eta + \tau_i + \varepsilon_{ij} = \mu_i + \varepsilon_{ij}$)

$\odot \quad \sum_{i=1}^{k} \tau_i = 0 \Rightarrow$ dummy variables: sum coding $\leftarrow$ LM, LNp. 8-11~16

When $n_1 = \cdots = n_k$ the sum codings are orthogonal to $\mathbb{1}$.

$\tau_1 + \tau_2 + \tau_3 + \tau_4 = 0 \Rightarrow \tau_4 = -(\tau_1 + \tau_2 + \tau_3) = \mu_4 - \bar{\mu}$

$$\begin{bmatrix}\mu_1\\\mu_2\\\mu_3\\\mu_4\end{bmatrix} = \eta\begin{bmatrix}1\\1\\1\\1\end{bmatrix} + \tau_1\begin{bmatrix}1\\0\\0\\0\end{bmatrix} + \tau_2\begin{bmatrix}0\\1\\0\\0\end{bmatrix} + \tau_3\begin{bmatrix}0\\0\\1\\0\end{bmatrix} + \tau_w\begin{bmatrix}0\\0\\0\\1\end{bmatrix} \cdots (*)$$

$$= \eta\begin{bmatrix}1\\1\\1\\1\end{bmatrix} + \tau_1\begin{bmatrix}1\\0\\0\\-1\end{bmatrix} + \tau_2\begin{bmatrix}0\\1\\0\\-1\end{bmatrix} + \tau_3\begin{bmatrix}0\\0\\1\\-1\end{bmatrix}$$

$\eta = \bar{\mu} = (\mu_1 + \cdots + \mu_4)/4$
$\tau_1 = \mu_1 - \bar{\mu}$
$\tau_2 = \mu_2 - \bar{\mu}$
$\tau_3 = \mu_3 - \bar{\mu}$

Interpretation of $\tau_i$'s, $\eta$

Recall. ANOVA decomposition in LNp.3-5 use the concept of sum codings. $\because$

- $\tau_1 = 0 \qquad \Rightarrow$ dummy variables: treatment coding

$$(*) = \eta\begin{bmatrix}1\\1\\1\\1\end{bmatrix} + \tau_2\begin{bmatrix}0\\1\\0\\0\end{bmatrix} + \tau_3\begin{bmatrix}0\\0\\1\\0\end{bmatrix} + \tau_w\begin{bmatrix}0\\0\\0\\1\end{bmatrix}$$

$\eta = \mu_1, \qquad \tau_2 = \mu_2 - \mu_1$
$\tau_3 = \mu_3 - \mu_1, \ \tau_4 = \mu_4 - \mu_1$

- $\eta = 0 \ (*) =$

$$\tau_1\begin{bmatrix}1\\0\\0\\0\end{bmatrix} + \tau_2\begin{bmatrix}0\\1\\0\\0\end{bmatrix} + \tau_3\begin{bmatrix}0\\0\\1\\0\end{bmatrix} + \tau_w\begin{bmatrix}0\\0\\0\\1\end{bmatrix} \quad \boxed{\tau_i = \mu_i}$$

$\Omega \ominus \omega = \text{span}\{\text{sum codings}\}$　ANOVA $\Rightarrow$ $H_0: \tau_1 = \cdots = \tau_k$　after the first 2 constraints

❖ **Reading**: textbook, 2.1 ← check LNp.2-34　　$H_0: \tau_1 = \cdots = \tau_k = 0$ ← on $\tau_i$'s are added
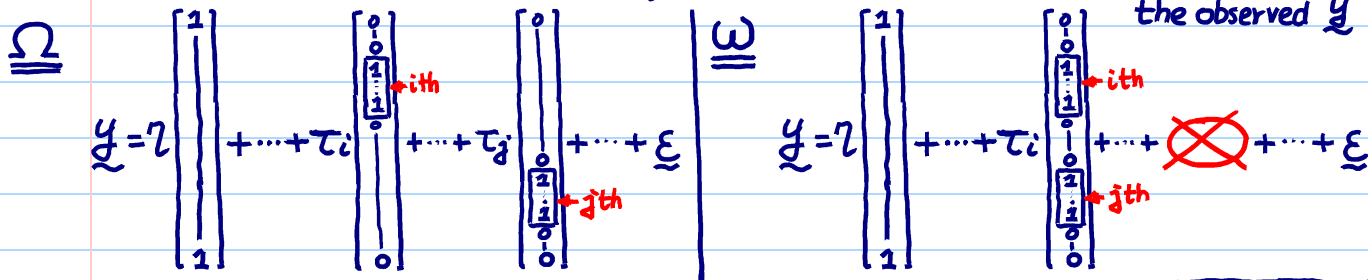
---

# Multiple Comparisons

answer "how different" problem: $\mu_i \overset{?}{=} \mu_j, \ \forall i < j$

- Consider the full model $y_{ij} = \eta + \tau_i + \varepsilon_{ij}$. For one pair, say $(i, j)$, of treatments, test $H_0^{ij}: \tau_i = \tau_j$ against $H_A^{ij}: \tau_i \neq \tau_j$.

decided ┬ before exp't
　　　　└ not related to the observed $y$

$\Longleftrightarrow \mu_i = \mu_j \ (\mu_i - \mu_j = 0)$

$\Omega$

$$\underset{\sim}{y} = \eta\begin{bmatrix}1\\ \vdots \\1\end{bmatrix} + \cdots + \tau_i\begin{bmatrix}0\\ \vdots\\1\\1\\ \vdots\\0\end{bmatrix}_{\leftarrow i\text{th}} + \cdots + \tau_j\begin{bmatrix}0\\ \vdots\\1\\ \vdots\\1\\0\end{bmatrix}_{\leftarrow j\text{th}} + \cdots + \underset{\sim}{\varepsilon}$$

$\omega$

$$\underset{\sim}{y} = \eta\begin{bmatrix}1\\ \vdots\\1\end{bmatrix} + \cdots + \tau_i\begin{bmatrix}0\\ \vdots\\1\\ \vdots\\1\\0\end{bmatrix}_{\leftarrow i\text{th}}^{\leftarrow j\text{th}} + \cdots + \bigotimes + \cdots + \underset{\sim}{\varepsilon}$$

$\Omega \ominus \omega = \text{span}\{\underset{\sim}{v}\}$, where



residual space

How can we derive it ?

$$\underset{\sim}{v} = \begin{bmatrix}\vdots\\0\\+1/n_i\\ \vdots\\+1/n_i\\0\\ \vdots\\-1/n_j\\ \vdots\\-1/n_j\\0\\ \vdots\end{bmatrix}_{\leftarrow j\text{th}}^{\leftarrow i\text{th}}$$

$\Omega \ominus \omega \ (\dim = 1)$　$\hat{\varepsilon}_\Omega = \underset{\sim}{y} - \bar{y}_{i\cdot}$

$$F_{ij} = \frac{RSS_\omega - RSS_\Omega /1}{RSS_\Omega /df_\Omega} \sim F_{1, df_\Omega} \ (\text{under} \ \omega)$$

$\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega = P_{\Omega \ominus \omega}\underset{\sim}{y}$

$$= \frac{\underset{\sim}{y}^T \underset{\sim}{v}}{\|\underset{\sim}{v}\|^2}\underset{\sim}{v} = \frac{\bar{y}_{i\cdot} - \bar{y}_{j\cdot}}{\sqrt{\frac{1}{n_i} + \frac{1}{n_j}}} \cdot \frac{\underset{\sim}{v}}{\|\underset{\sim}{v}\|}$$

$\Omega \ (\dim = k)$　$\omega \ (\dim = k-1)$

$RSS_\omega - RSS_\Omega = \|\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega\|^2$
$RSS_\Omega = \|\hat{\varepsilon}_\Omega\|^2$

jointly made by Jeff Wu (GT, USA) and S.-W. Cheng (NTHU, Taiwan)

- It is common to use the _t_-test and the _t_-statistic

$t_{ij}^2 = F_{ij}$

$$t_{ij} = \frac{\bar{y}_{i\cdot} - \bar{y}_{j\cdot} - 0}{\hat{\sigma}\sqrt{1/n_i + 1/n_j}},$$

$\mu_i \quad \mu_j$

$\boxed{\bar{y}_{i\cdot}} \longleftrightarrow \boxed{\bar{y}_{j\cdot}}$

$Var(\bar{y}_{i\cdot} - \bar{y}_{j\cdot})$
$= Var(\bar{y}_{i\cdot}) + Var(\bar{y}_{j\cdot})$
$= \sigma^2/n_i + \sigma^2/n_j$

where $n_i$ = number of observations for treatment $i$,

$\hat{\sigma}^2 = RSS_\Omega / df_\Omega$ in ANOVA; declare "treatments

$i$ and $j$ different at level α" if

$F_{1,df_\Omega} \to$ null dist.

$df_\Omega$

$|t_{ij}| > t_{N-k, \frac{\alpha}{2}}.$

null of ANOVA in LNp.3-b

multiple testing

$R_{12} \quad R_{13}$
$R_{23}$

$R_{ij} \equiv$ rejection region of $H_0^{ij}$
$P(R_{ij} \mid \mu_i = \mu_j \,(\tau_i = \tau_j)) = \alpha$

EER
$= P(\cup R_{ij} \mid \mu_1 = \cdots = \mu_k)$
$> \alpha \leftarrow$ why not good?
$\hookleftarrow$ usually

- Suppose $k'$ tests are performed to test $H_0 : \tau_1 = \cdots = \tau_k$.

e.g. $\binom{k}{2}$

_Experimentwise error rate_ (EER) = Probability of declaring at least one pair of treatments significantly different under $H_0$. Need to use multiple comparisons to control EER.

$\bigcap_{i,j} H_0^{ij}$ $\leftarrow$ union-intersection test

$H_1$

$\bigcup_{i,j} H_1^{ij}$

$t_{16, \frac{0.05}{2}} = 2.12$

$t_{ij}$

| A vs. B | A vs. C | A vs. D | B vs. C | B vs. D | C vs. D |
|---------|---------|---------|---------|---------|---------|
| −0.87   | 1.85    | 2.14    | 2.72    | 3.01    | 0.29    |

DC    A B

Note. deductive logic does not hold (Why?)

$|\bar{y}_{i\cdot} - \bar{y}_{j\cdot}|$
$> \hat{\sigma}\sqrt{\frac{1}{n_i} + \frac{1}{n_j}} \cdot t_{N-k, \frac{\alpha}{2}}$

a scale

$\bar{y}_{(1)\cdot} \; \bar{y}_{(2)\cdot} \cdots \bar{y}_{(k-1)\cdot} \; \bar{y}_{(k)\cdot}$

$\mu_{(1)} \; \mu_{(2)} \cdots \mu_{(k-1)} \; \mu_{(k)}$

$\mu_{(1)}, \cdots, \mu_{(k)}$ : sorted $\mu_1, \cdots, \mu_k$
$(\mu_{(1)} \leq \mu_{(2)} \leq \cdots \leq \mu_{(k)})$

$\bar{y}_{(1)\cdot}, \cdots, \bar{y}_{(k)\cdot}$ : order statistics of $\bar{y}_{1\cdot}, \cdots, \bar{y}_{k\cdot}$.

---

# Bonferroni Method

c.f. $\boxed{\frac{\alpha}{2}}$

◉ Declare "$\tau_i$ different from $\tau_j$ at level α" if $|t_{ij}| > t_{N-k, \frac{\alpha}{2k'}}$, where $k'$ = number of tests.

$\boxed{EER \leq \alpha}$

very conservative when $k'$ is large

- For one-way layout with $k$ treatments, $k' = \binom{k}{2} = \frac{1}{2}k(k-1)$, as $k$ increases, $k'$ increases, and the critical value $t_{N-k, \frac{\alpha}{2k'}}$ gets bigger (i.e., method less powerful in detecting differences).
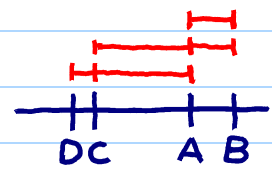
some indep. assumptions btwn the test stat. of these tests, e.g., indep. $\bar{y}_{i\cdot}$'s

- Advantage: It works without requiring independence assumption.

can be widely applied.

- For pulp experiment, take α = 0.05, $k$ = 4, $k'$ = 6, $t_{16, 0.05/12}$ = 3.008. Among the 6 $t_{ij}$-values (see LNp.3-10), only the $t$-value for B-vs-D, 3.01, is larger. Declare "B and D different at level 0.05".

$2 \times 6$

$P(\cup R_{ij} \mid \cap H_0^{ij})$
$\leq \sum_{i,j} P(R_{ij} \mid H_0^{ij}) = k' \cdot \alpha' = \alpha$

$\Rightarrow \alpha' = P(R_{ij} \mid H_0^{ij}) = \alpha/k'$

DC    A B

For simplicity, assume $n_1 = \cdots = n_k = n$

## Tukey Method
← Same procedure can be applied to unequal sample size case → Tukey-Kramer test

- Declare "$\tau_i$ different from $\tau_j$ at level $\alpha$" if

$$\because \sqrt{1/n + 1/n} = \sqrt{2}/\sqrt{n}$$

$$|t_{ij}| > \frac{1}{\sqrt{2}} q_{k,N-k,\alpha}, \qquad \boxed{EER \le \alpha}$$

(LNp.3-10) where $q_{k,N-k,\alpha}$ is the upper $\alpha$-quantile of the **Studentized range** (**SR**) distribution with parameter $k$ and $N-k$ degrees of freedom. (see distribution table on LNp.3-13)

- For pulp experiment,

LNp.3-13 → $\sqrt{RSS_n/df_n}$ (LNp.3-9)

$$\frac{1}{\sqrt{2}} q_{k,N-k,0.05} = \frac{1}{\sqrt{2}} q_{4,16,0.05} = \frac{4.05}{\sqrt{2}} = 2.86.$$

Again only B-vs-D has larger $t_{ij}$-value than 2.86 (See LNp.3-10). Tukey method is more powerful than Bonferroni method because 2.86 is smaller than 3.01 (why?)

Q: Among the 3 critical values of $|t_{ij}|$, which one is the smallest? the largest?

---

- compare $\bar{y}_{i\cdot} - \bar{y}_{j\cdot}$, $\forall(i,j)$    p. 3-12
- For $1 \le i < j \le k$,    check LNp.3-10
  $|\bar{y}_{i\cdot} - \bar{y}_{j\cdot}| \le \bar{y}_{(k)\cdot} - \bar{y}_{(1)\cdot}$
- $P(\cup R_{ij} | \cap H_o^{ij})$ ← EER
  $= P(\text{at least one } |\bar{y}_{i\cdot} - \bar{y}_{j\cdot}| \text{ larger than } \underline{c^*} | \cap H_o^{ij})$
  $= P(\bar{y}_{(k)\cdot} - \bar{y}_{(1)\cdot} > c^* | \cap H_o^{ij})$
  $= \alpha$
- Under $\cap H_o^{ij}$, $\bar{y}_{1\cdot}, \cdots, \bar{y}_{k\cdot}$
  indep. $N(\mu, \sigma^2/n) \Rightarrow$
  $\dfrac{\sqrt{n}\,(\bar{y}_{(k)\cdot} - \bar{y}_{(1)\cdot})}{\hat{\sigma}} \sim SR_{k,\nu}$
  where $\nu\hat{\sigma}^2 \sim \sigma^2 \chi_\nu^2$ and indep. of $\bar{y}_{1\cdot}, \cdots, \bar{y}_{k\cdot}$.
  $\Rightarrow c^* = \hat{\sigma}/\sqrt{n} \cdot SR_{k,\nu,\alpha}$
  $\Rightarrow \sqrt{n}\,|\bar{y}_i - \bar{y}_j|/\hat{\sigma} > SR_{k,\nu,\alpha}$

Q: For the case in LNp.3-10, after getting $\bar{y}_{i\cdot}$'s, want to test $H_o^{ij}: \mu_D = \mu_B$ (only one test). Which critical value should we use? $H_o: \mu_{(1)} = \mu_{(k)}$?

---

## Selected values of $q_{k,\nu,\alpha}$ for $\alpha = 0.05$
$\nu = N-k$

increasing (why?)

| $\nu$ | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 17.97 | 26.98 | 32.82 | 37.08 | 40.41 | 43.12 | 45.40 | 47.36 | 49.07 | 50.59 | 51.96 | 53.20 | 54.33 | 55.36 |
| 2 | 6.08 | 8.33 | 9.80 | 10.88 | 11.74 | 12.44 | 13.03 | 13.54 | 13.99 | 14.39 | 14.75 | 15.08 | 15.38 | 15.65 |
| 3 | 4.50 | 5.91 | 6.82 | 7.50 | 8.04 | 8.48 | 8.85 | 9.18 | 9.46 | 9.72 | 9.95 | 10.15 | 10.35 | 10.52 |
| 4 | 3.93 | 5.04 | 5.76 | 6.29 | 6.71 | 7.05 | 7.35 | 7.60 | 7.83 | 8.03 | 8.21 | 8.37 | 8.52 | 8.66 |
| 5 | 3.64 | 4.60 | 5.22 | 5.67 | 6.03 | 6.33 | 6.58 | 6.80 | 6.99 | 7.17 | 7.32 | 7.47 | 7.60 | 7.72 |
| 6 | 3.46 | 4.34 | 4.90 | 5.30 | 5.63 | 5.90 | 6.12 | 6.32 | 6.49 | 6.65 | 6.79 | 6.92 | 7.03 | 7.14 |
| 7 | 3.34 | 4.16 | 4.68 | 5.06 | 5.36 | 5.61 | 5.82 | 6.00 | 6.16 | 6.30 | 6.43 | 6.55 | 6.66 | 6.76 |
| 8 | 3.26 | 4.04 | 4.53 | 4.89 | 5.17 | 5.40 | 5.60 | 5.77 | 5.92 | 6.05 | 6.18 | 6.29 | 6.39 | 6.48 |
| 9 | 3.20 | 3.95 | 4.41 | 4.76 | 5.02 | 5.24 | 5.43 | 5.59 | 5.74 | 5.87 | 5.98 | 6.09 | 6.19 | 6.28 |
| 10 | 3.15 | 3.88 | 4.33 | 4.65 | 4.91 | 5.12 | 5.30 | 5.46 | 5.60 | 5.72 | 5.83 | 5.93 | 6.03 | 6.11 |
| 11 | 3.11 | 3.82 | 4.26 | 4.57 | 4.82 | 5.03 | 5.20 | 5.35 | 5.49 | 5.61 | 5.71 | 5.81 | 5.90 | 5.98 |
| 12 | 3.08 | 3.77 | 4.20 | 4.51 | 4.75 | 4.95 | 5.12 | 5.27 | 5.39 | 5.51 | 5.61 | 5.71 | 5.80 | 5.88 |
| 13 | 3.06 | 3.73 | 4.15 | 4.45 | 4.69 | 4.88 | 5.05 | 5.19 | 5.32 | 5.43 | 5.53 | 5.63 | 5.71 | 5.79 |
| 14 | 3.03 | 3.70 | 4.11 | 4.41 | 4.64 | 4.83 | 4.99 | 5.13 | 5.25 | 5.36 | 5.46 | 5.55 | 5.64 | 5.71 |
| 15 | 3.01 | 3.67 | 4.08 | 4.37 | 4.59 | 4.78 | 4.94 | 5.08 | 5.20 | 5.31 | 5.40 | 5.49 | 5.57 | 5.65 |
| 16 | 3.00 | 3.65 | 4.05 | 4.33 | 4.56 | 4.74 | 4.90 | 5.03 | 5.15 | 5.26 | 5.35 | 5.44 | 5.52 | 5.59 |

$\alpha$=upper tail probability, $\nu$=degrees of freedom, $k$=number of treatments

decreasing (Why?)

For complete tables corresponding to various values of $\alpha$ refer to Appendix E.

❖ **Reading**: textbook, 2.2

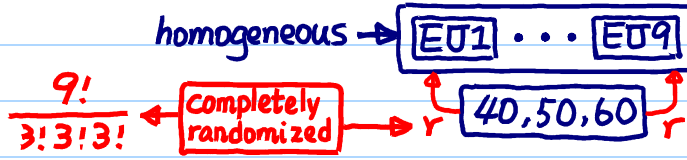# One-Way ANOVA with a Quantitative Factor

cf. Qualitative

✿ response : bonding strength

♣ (treament) factor : power (quantitative)

3 levels — 40, 50, 60

[equally spaced]

✿ Exp'tal units : one composite

9 EUs

✿ Each treatment repeats 3 times ( 3 replicates)

● **Data** :

Design matrix
power strength

$$
\begin{bmatrix}
40 \rightarrow 25.66 \\
40 \rightarrow 28.00 \\
40 \rightarrow 20.65 \\
50 \\
50 \\
50 \\
60 \\
60 \\
60 \rightarrow 35.66
\end{bmatrix}
$$

$y$ = bonding strength of composite material,

$x$ = laser power at 40, 50, 60 watt.

homogeneous → | EU1 | ··· | EU9 |

9!
3! 3! 3! ← [Completely randomized] → r | 40, 50, 60 | r

Table 2: Strength Data, Composite Experiment

Initial Data Analysis : scatter plot

the line is meaningful only when the factor is quantitative

major difference between quantitative & qualitative factors

| Laser Power (watts) | | |
|---|---|---|
| 40 | 50 | 60 |
| 25.66 | 29.15 | 35.73 |
| 28.00 | 35.09 | 39.56 |
| 20.65 | 29.79 | 35.66 |

conceptual model

$\mu_x = E(y_x)$

$y_x = \mu_x + \varepsilon$

$\mu_x = \beta_0 + \beta_1 x$

✿ $\mu_x = \beta_0 + \beta_1 x + \beta_2 x^2$

$\mu_x = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$

$\mu_x = \beta_0 + \beta_1 x \log x + \beta_2 e^x$

⋮