LM, LNp.7-1~2 → **Residual Analysis: Theory**

- Theory: define the **residual** for the $i^{th}$ observation $(x_i, y_i)$ as

$$\widehat{\varepsilon}_i = r_i = y_i - \hat{y}_i, \qquad \hat{y}_i = \mathbf{x}_i^T \hat{\beta},$$  ⎫ $i$th row of model matrix X

$\widehat{\varepsilon}_i$ often used to check assumptions:
① $\varepsilon_i$'s $\overset{i.i.d}{\sim}$ $N(0,\sigma^2)$
② $E(Y) = X\beta$ is a correct mean structure

$\hat{y}_i$ contains <u>information given by the</u> <u>model</u>; $r_i$ is the "<u>difference</u>" between $y_i$ (<u>observed</u>) and $\hat{y}_i$ (<u>fitted</u>) and contains <u>information on possible</u> *model inadequacy*.

- <u>Vector of residuals</u> $\mathbf{r} = (r_1, \ldots, r_N)^T = \mathbf{y} - \mathbf{X}\hat{\beta}. = (I - H)Y$
  ⎿ hat matrix

  overall pattern  → individual observation (outlier)

- Under the model assumption $E(\mathbf{y}) = \mathbf{X}\beta$, it can be shown that

  (a) $E(\mathbf{r}) = \underline{0}$,   ∵ $\underline{r} \perp \hat{\underline{y}}$   ⎿ i.e., $X\beta$ is <u>correct model</u>

  (b) $\mathbf{r}$ and $\hat{\mathbf{y}}$ are <u>independent</u>,

  Note. Variance of $r_i$ ∝ $1 - h_i$ (leverage) $= 1 - H_{ii}$ ↳ ∝ Mahananobis dist. btwn design pts & $\bar{X}$

  cov($\underline{r}$) $= \sigma^2(I-H)$

  (c) <u>variances of $r_i$ are nearly constant</u> for "<u>nearly balanced</u>" designs.

① $\underline{Y} = \underline{X\beta} + \underline{\varepsilon} = \underline{\hat{Y}} + \underline{\hat{\varepsilon}}$

② $Y = \underline{X_1\beta_1} + \underline{X_2\beta_2} + \varepsilon = (X_1\beta_1 + \underline{H_1X_2\beta_2}) + (\underline{(I-H_1)X_2\beta_2} + \varepsilon) = \underline{\hat{Y}_{X_1}} + \underline{\hat{\varepsilon}_{X_1}})$

⎿ Under fitted model: $Y = X_1\underline{\beta}_1 + \varepsilon^*$ ($H_1 = X_1(X_1^TX_1)^{-1}X_1^T$)

---

$\sigma_x \leftrightarrow \mu_x$   **Residual Plots** ← LM, LNp.7-7~9

- Plot $r_i$ vs. $\hat{y}_i$ (see Figure 1): <u>It should appear as a parallel band around 0</u>. ⎿ LNp.24
  <u>Otherwise</u>, it would suggest <u>model violation</u>. If <u>spread of $r_i$ increases</u> as $\hat{y}_i$ <u>increases, error variance of $y$ increases with mean of $y$</u>. May need a <u>transformation of $y$</u>. (Will be explained in <u>future lecture</u>.)

  ⌐ LNp.25

- ⊙ Plot $r_i$ from <u>replicates per treatment</u> (see Figure 2): to see  ← $\hat{\underline{r}}$ vs. predictor
  <u>if error variance depends on treatment</u>.  Note. saturated model

  qualitative factor → **Box plots**

- ⊙ Plot $r_i$ vs. $x_i$: If <u>not a parallel band around 0</u>, relationship between
  quantitative factor  <u>$y_i$ and $x_i$ not fully captured</u>, revise the $\mathbf{X}\beta$ part of the model.

- Plot $r_i$ vs. <u>time sequence</u>: to see if there is a  if available (or run order, measure order, ...)
  <u>time trend</u> or <u>autocorrelation over time</u>.

null plot

# Plot of $r_i$ vs. $\hat{y}_i$

no replicate data

simple
model            $\hat{\varepsilon}$

complex
model            $\hat{\varepsilon}$

saturated
model            $\hat{\varepsilon}$

replicates

allow "pure
error" variance
estimation



$\bar{y}_D$   $\bar{y}_C$   $\bar{y}_A$   $\bar{y}_B$

fitted

$\widehat{y}_i \leftarrow$ a function of treatments A,B,C,D

Figure 1: $r_i$ vs. $\hat{y}_i$, Pulp Experiment
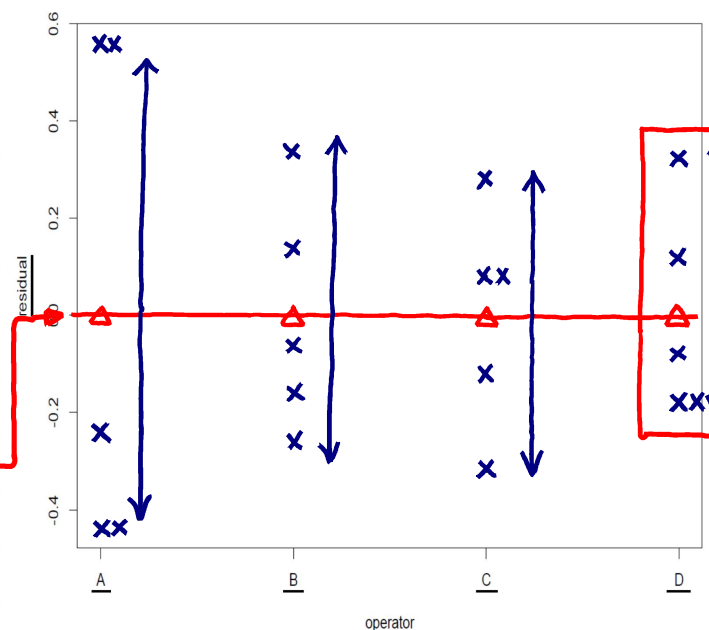
$\llcorner$ LNp.1

# Plot of $r_i$ vs. treatment

fitted model:
saturated

Same 4 groups
of residuals as
in the residual
plot of $r_i$ vs.
$\widehat{y}_i$ (LNp.3-24)

no need to
check the
trend in the
mean of
residuals



$\sum_j \widehat{\varepsilon}_{ij} = 0$
for any $i$

$$X = \begin{bmatrix} 1 & 1 & 0 & & 0 \\ 1 & 1 & 0 & & 0 \\ 1 & 1 & 0 & & \\ 1 & 1 & 0 & & \\ 1 & 1 & 1 & \cdot & \\ 1 & 1 & 0 & & \\ & 0 & 1 & & \\ 1 & 0 & 1 & & 1 \\ 1 & 0 & 0 & & 1 \\ & & & & 1 \end{bmatrix}$$

Figure 2: $r_i$ vs. treatment, Pulp Experiment

$\llcorner$ LNp.1

$Y_1, \ldots, Y_n$ iid ~ cdf $\underline{\underline{F}}$ — ## Box-(Whisker) Plot

- A powerful graphical display (due to Tukey) to capture the location, dispersion, skewness and extremity of a distribution. See Figure 3.

  ↳ LNp 27

- $Q_1$ = lower quartile (25% quantile), $Q_3$ = upper quartile (75% quantile), $Q_2$ = median (50% quantile, estimate of *location* parameter) is the white line in the box. $Q_1$ and $Q_3$ are boundaries of the *black box*.

- *IQR* = interquartile range (length of box) = $Q_3 - Q_1$: measure of *dispersion*.

- Minimum and maximum of **observed** values within
$$[Q_1 - 1.5 \times IQR, \ \ Q_3 + 1.5 \times IQR]$$
are denoted by two *whiskers*. Any values outside the whiskers are regarded as extreme values and displayed (possible *outliers*).

- If $Q_1$ and $Q_3$ are not symmetric around the median, it indicates *skewness*.

- Side-by-side box plots (LNp. 3-2~3) are useful to compare the difference between the distributions of several groups of data.
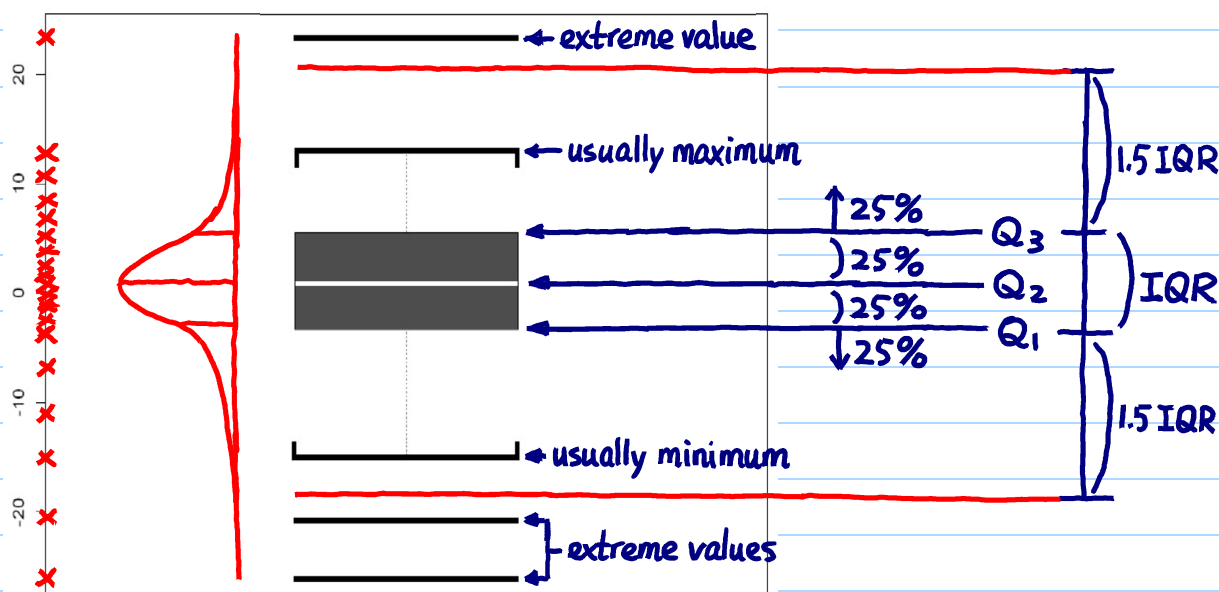
---

p. 3-27

# Box-(Whisker) Plot



Figure 3: Box-Whisker Plot

# Normal Probability Plot ← Q-Q plot (LM, LNp.7-15~16)

⊙ Original purpose : To test if a distribution is normal, e.g., if the residuals
follow a normal distribution (see Figure 5). **Q: Why need normality for error?**
→ can be used to identify outlier ↳LNp.30

⊙ More powerful use in factorial experiments (discussed in Units 5 and 6).
→ used to identify significant effects $\beta_i$'s ←cf.→ replace t-tests

• Let $r_{(1)} \leq \ldots \leq r_{(N)}$ be the ordered residuals. The cumulative
probability for $r_{(i)}$ is $p_i = (i-0.5)/N$. Thus the plot of ← LNp.29
$p_i$ vs. $r_{(i)}$ should be S-shaped as in Figure 4(a) if the errors
are normal. By transforming the scale of the horizontal axis,
the S-shaped curve is straightened to be a line (see Figure 4(b)).
↳ LNp.29

true cdf ←
empirical cdf ↘

• **Normal probability plot** of residuals :

$$\left( \Phi^{-1}\left(\frac{i-0.5}{N}\right), r_{(i)} \right), \qquad i = 1,\ldots,N, \qquad \Phi = \text{normal cdf.}$$

If the errors are normal, it should plot roughly as a straight line. See
Figure 5.
↳ LNp. 30

---

# Regular and Normal Probability Plots of Normal

(a)   **Φ: cdf of N(0,1)**     **empirical cdf of $\hat{\mathcal{E}}_i$**     CDF     (b)   $\hat{\mathcal{E}}_{(i)}$     **y=x**

$\frac{i-0.5}{N}$     ←$\frac{1}{N}$

$\Phi^{-1}(\frac{i-0.5}{N})$     ↑└$\hat{\mathcal{E}}_{(i)}$     $\Phi^{-1}$
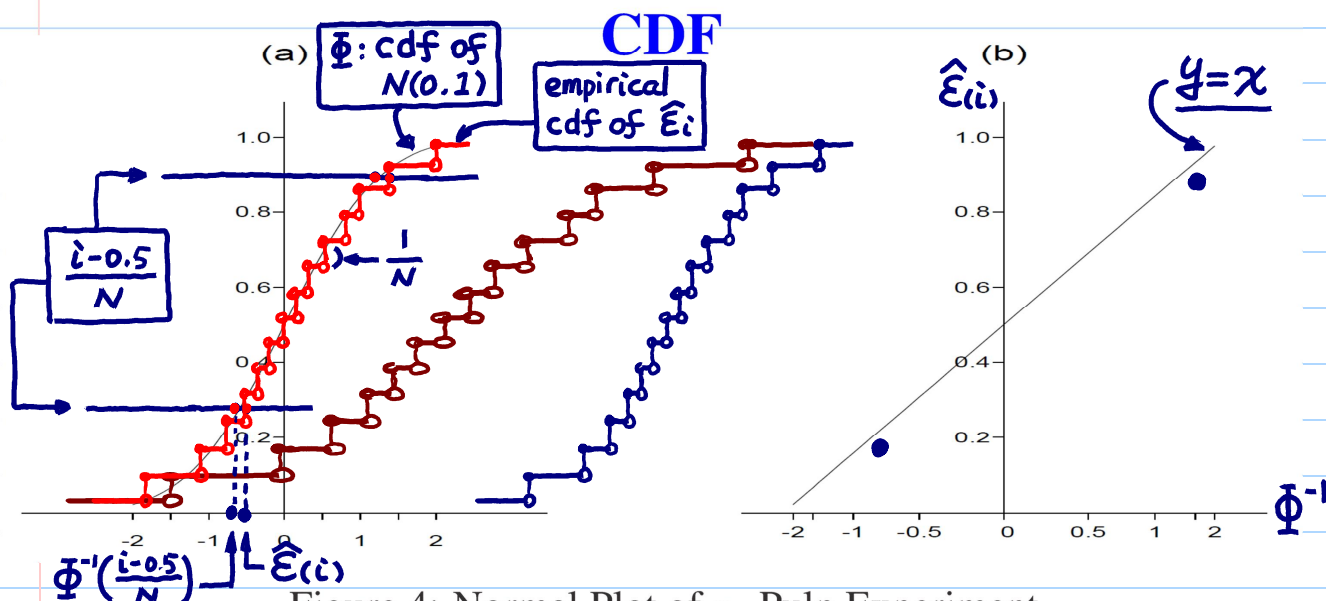
Figure 4: Normal Plot of $r_i$, Pulp Experiment

$Z_1, \ldots, Z_N$ iid $\sim N(\mu, \sigma^2)$

$Z_i = \sigma W_i + \mu \Leftarrow W_i = (Z_i - \mu)/\sigma$ iid $\sim N(0,1)$

Normal probability plot of $W_i$'s ⇒ $W_{(i)}$ vs. $\Phi^{-1}$: $y = x$

Normal probability plot of $Z_i$'s ⇒ $Z_{(i)}$ vs. $\Phi^{-1}$: $y = \sigma x + \mu$
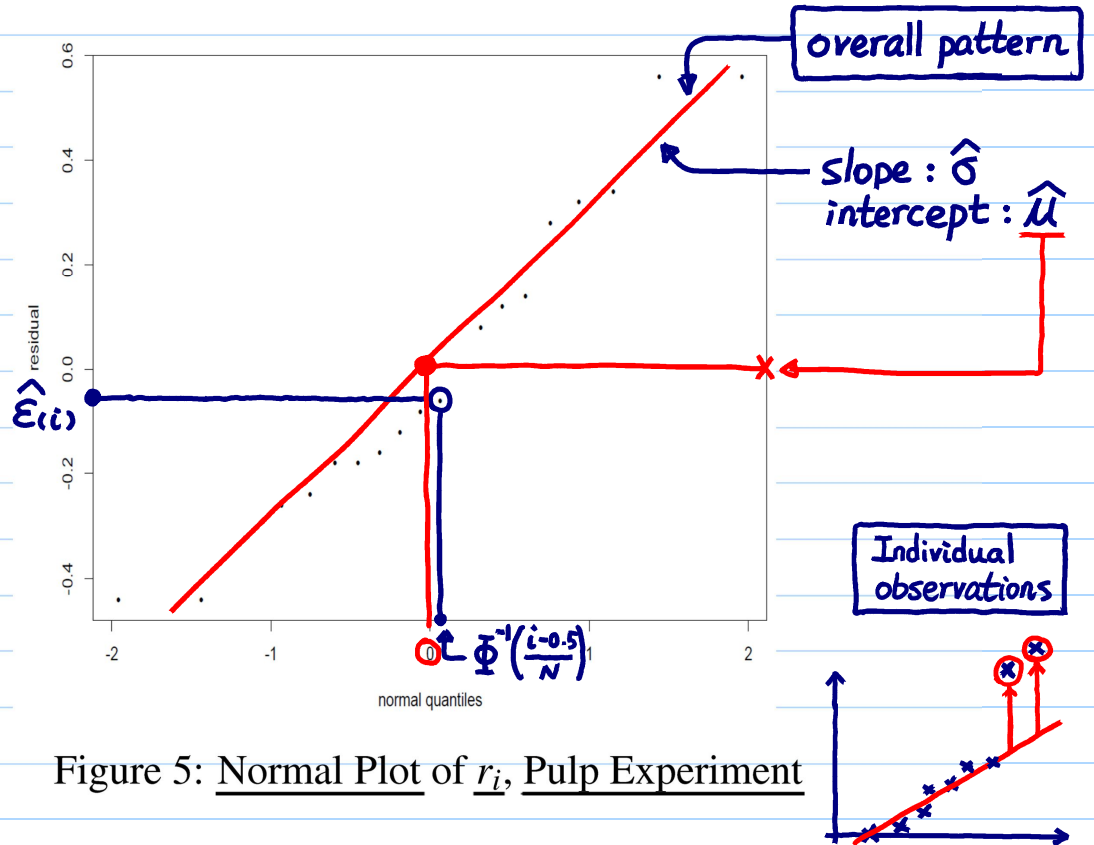
# Normal Probability Plot : Pulp Experiment

overall pattern

slope : $\hat{\sigma}$
intercept : $\hat{\mu}$

$\hat{\varepsilon}_{(i)}$ residual

$\Phi^{-1}\left(\frac{i-0.5}{N}\right)$

normal quantiles

Individual observations

Figure 5: Normal Plot of $r_i$, Pulp Experiment

❖ **Reading**: textbook, 2.6

---

# Pulp Experiment Revisited

Q: Are A.B.C.D a representitive sample of the population ?

population

LNp.1

- Compare the 2 scenarios

LNp.4
(S1) plant has only 4 operators (or only interested in these 4 operators)

effects of operators

after sampling & coditioning
— $\tau_i$'s: parameters (unknown fixed values)
— interest: difference btwn the 4 specific $\tau_i$'s

cf.　(S2) 4 operators randomly sampled from a large population of operators

before sampling
— $\tau_i$'s: random variables

$\mu_i = E(y_{ij})$
— interest: difference btwn all operators in this population

$\tau_i$'s

dist. of $\tau_i$'s

1 ② 3 ... ↑...↑...↑...R operators
　　B　　　C　A　　D

- In the pulp experiment the effects $\tau_i$ are called *fixed* effects because the interest was in comparing the four *specific* operators in the study. If these four operators were chosen randomly from the population of operators in the plant, the interest would usually be in the variation among all operators in the population. Because the observed data are from operators randomly selected from the population, the variation among operators in the *population* is referred to as *random* effects.

conditioned on these 4 operators

p. 3-32

**fixed effect model**

- **One-way random effects model (REM $\xleftrightarrow{cf.}$ FEM) :**

FEM in LNp.3-4 & 3-8 $\xleftrightarrow{cf.}$

intercept, parameter

$\xrightarrow{cf.}$ whole- & sub-plot errors in LNp. 4-47~48

$$y_{ij} = \eta + \tau_i + \varepsilon_{ij},$$

r.v.

$y_{ij} \sim N(\eta, \sigma_\tau^2 + \sigma^2)$

Q: $y_{ij}$'s indep.?

where $\varepsilon_{ij}$'s:　independent error terms with $N(0, \sigma^2)$,

$\tau_i$'s:　independent $N(0, \sigma_\tau^2)$,

why?

and $\tau_i$ and $\varepsilon_{ij}$ are independent (Why? Give an example.);

REM: parameters $\eta, \sigma_\tau^2, \sigma^2$

$\sigma^2$ and $\sigma_\tau^2$ are the two *variance components* of the model. check $\Sigma^*$

The variance among operators in the population is measured by $\sigma_\tau^2$.

$Y = \begin{bmatrix} y_{11} \\ \vdots \\ y_{1n_1} \\ y_{21} \\ \vdots \\ y_{2n_2} \\ \vdots \\ y_{k1} \\ \vdots \\ y_{kn_k} \end{bmatrix}$

**REM.** (1) $E(Y) = \eta \mathbf{1}$

(2) $cov(y_{11}, y_{12}) = Cov(\eta + \tau_1 + \varepsilon_{11}, \eta + \tau_1 + \varepsilon_{12})$
$= \sigma_\tau^2$

$cov(y_{11}, y_{21}) = Cov(\eta + \tau_1 + \varepsilon_{11}, \eta + \tau_2 + \varepsilon_{22})$
$= 0$

$cov(Y) = \underset{\underset{\Sigma^*}{=\!=\!=}}{\begin{bmatrix} \boxed{}_{n_1 \times n_1} & & \underset{\sim}{0} \\ & \boxed{}_{n_2 \times n_2} & \\ \underset{\sim}{0} & & \boxed{} \end{bmatrix}} \xrightarrow{} \begin{bmatrix} \sigma_\tau^2 + \sigma^2 & & \sigma_\tau^2 \\ & \ddots & \\ \sigma_\tau^2 & & \sigma_\tau^2 + \sigma^2 \end{bmatrix}$

$= \sigma_\tau^2 \begin{bmatrix} 1 & & 1 \\ 1 & \ddots & 1 \\ 1 & & 1 \end{bmatrix} + \sigma^2 \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix}$ $\boxed{\mathcal{X}}$

**FEM.** $y_{ij} \xrightarrow{indep.} N(\eta + \tau_i, \sigma^2)$

(1) $E(Y) = X\underline{\beta} \leftarrow \eta + \tau_i$

(2) $cov(Y) = \sigma^2 I$ $\boxed{\mathcal{X}}$

★ relationship btwn $y_x$ & $\mathcal{X}$
- In FEM, $E(Y)$
- In REM, $cov(Y)$

---

p. 3-33

$H_0: \sigma_\tau^2 = 0$ in LNp 4-63 $\xleftrightarrow{cf.}$

# One-way Random Effects Model: ANOVA

balance ● In the following, assume $n_1 = \cdots = n_k = n$.

● The null hypothesis for the FEM:

$N(\eta\mathbf{1}, \Sigma^*) \sim Y$

check LNp.5

$\Omega \ominus \omega$

same null dist.

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

should be replaced by　meaning?

$P_0, P_1, P_2$: projection matrix onto $\omega$, $\Omega \ominus \omega$, $\Omega^\perp$

$Y \sim N(\eta\mathbf{1}, \sigma^2 I) \to$ $H_0^* : \sigma_\tau^2 = 0.$

Under $H_0^*$, the $F$-test and the

ANOVA table in LNp. 3-6 still holds.

$\Omega$ (dim = $k$)　span$\{\mathbf{1}\} = \omega$ (dim = 1)

● Reason: under $H_0^*$,

$E(P_0 Y) = P_0 \eta\mathbf{1} = \eta\mathbf{1}$

(under $H_0^*$) $SSTr \sim \sigma^2 \chi_{k-1}^2$, and

$E(P_1 Y) = E(P_2 Y) = \underset{\sim}{0}$

(under $H_0^*$ & $H_A^*$) $SSE \sim \sigma^2 \chi_{N-k}^2$,

$\dfrac{\dfrac{SSTr}{k-1}}{\dfrac{SSE}{N-k}}$ and they are independent.

Therefore the $F$-test has the

distribution $F_{k-1, N-k}$ under $H_0^*$.

$cov(P_1 Y, P_2 Y) = P_1 cov(Y) P_2^T$
$= P_1 \Sigma^* P_2^T = \underset{\sim}{0}$

columns of $P_1$ & $P_2$ are eigenvectors of $\Sigma^*$ & $P_1 P_2 = \underset{\sim}{0}$

$\Omega^\perp, \Omega$ : different eigenspace of $\Sigma^*$
$\sigma^2 \quad n\sigma_\tau^2 + \sigma^2$ : eigenvalue (exercise, use the vectors in LNp.3-4)

$SSTr = \|P_1 Y\|^2 = \sum_i n(\bar{y}_{i\cdot} - \bar{y}_{\cdot\cdot})^2$ (LNp.3-5)
$= \sum_i (\sqrt{n}\,\bar{y}_{i\cdot} - \sqrt{n}\,\bar{y}_{\cdot\cdot})^2 \sim (n\sigma_\tau^2 + \sigma^2)\chi_{k-1}^2$

(O)

$\sqrt{n}\,\bar{y}_{i\cdot} = \sqrt{n}(\eta + \tau_i + \bar{\varepsilon}_{i\cdot}) \overset{iid}{\sim} N(\sqrt{n}\,\eta, n\sigma_\tau^2 + \sigma^2)$
$\bar{y}_{\cdot\cdot} = \bar{\bar{y}}_{i\cdot}$ check LNp 36

$SSE = \|P_2 Y\|^2 = \sum_i \sum_j (y_{ij} - \bar{y}_{i\cdot})^2$ (LNp.3-5)

$y_{ij} = \eta + \tau_i + \varepsilon_{ij}$

$= \sum_j (\varepsilon_{ij} - \bar{\varepsilon}_{i\cdot})^2 \sim \sigma^2 \chi_{N-k}^2$
$\sim \sigma^2 \chi_{n-1}^2$ $\quad k(n-1)$

# ANOVA Tables ($n_i = n$)

- We can apply the *same* ANOVA and $F$-test in the <u>fixed effects</u> case for analyzing data.
  - same test statistic
  - same null dist

**different** ← cf. E(MS) in LNp.3-6

**ANOVA table (FEM) in LNp.3-6** cf.

| Source | d.f. | SS | MS | E(MS) ← Under $H_0^* \cup H_A^*$ |
|--------|------|-----|-----|------|
| treatment | $k-1$ | $SSTr$ | $MSTr = \frac{SSTr}{k-1}$ | $\sigma^2 + n\sigma_\tau^2$ ← Under $H_0^*$, $\sigma_\tau^2 = 0$ $E(MSTr) = \sigma^2$ |
| residual | $N-k$ | $SSE$ | $MSE = \frac{SSE}{N-k}$ | $\sigma^2$ |
| total | $N-1$ | | | |

**ANOVA result (FEM) in LNp. 3-7** cf.

Pulp Experiment

| Source | d.f. | SS | MS | E(MS) |
|--------|------|-----|-----|-------|
| treatment | 3 | 1.34 | 0.447 | $\sigma^2 + 5\sigma_\tau^2$ |
| residual | 16 | 1.70 | 0.106 | $\sigma^2$ |
| total | 19 | 3.04 | | |

- However, we need to <u>compute</u> the <u>expected mean squares</u> under the <u>alternative</u> of $\sigma_\tau^2 > 0$,

  (i) for <u>sample size</u> <u>determination</u>, and
  (ii) to <u>estimate</u> the <u>variance components</u>. ($\sigma_\tau^2$ & $\sigma^2$)

# Expected Mean Squares for Treatments

- <u>Equation (1)</u> <u>holds</u> <u>independent of</u> $\sigma_\tau^2$,
  (LNp.3-33) $\sigma^2 \chi_{N-k}^2 \sim$

  **SSE/N-k : an unbiased estimator of $\sigma^2$**

$$E(MSE) = E\left(\frac{SSE}{N-k}\right) = \sigma^2. \tag{1}$$

**SSE only contains information of error var. component $\sigma^2$**

- Under the <u>alternative</u>: $\sigma_\tau^2 > 0$, and for $n_i = n$,
  (LNp.3-33) $(n\sigma_\tau^2 + \sigma^2)\chi_{k-1}^2 \sim$

$$E(MSTr) = E\left(\frac{SSTr}{k-1}\right) = \sigma^2 + n\sigma_\tau^2. \tag{2}$$

$$E\left(\frac{\frac{SSTr}{k-1} - \frac{SSE}{N-k}}{n}\right) = \sigma_\tau^2$$

→ an unbiased estimator of $\sigma_\tau^2$

**SSTr contains information about factor var. component $\sigma_\tau^2$ error var. component $\sigma^2$**

(cf. E(SSTr) of FEM in LNp 3-6)

- For <u>unequal</u> $n_i$'s, <u>$n$</u> in (2) is <u>replaced by</u>

$$n' = \frac{1}{k-1}\left[\sum_{i=1}^{k} n_i - \frac{\sum_{i=1}^{k} n_i^2}{\sum_{i=1}^{k} n_i}\right].$$

(exercise) use (0) in LNp.3-33

$$y_{ij} = \eta + \tau_i + \varepsilon_{ij}$$

**Proof of (2)**

$z_1, \ldots, z_k \overset{\text{i.i.d.}}{\sim} N(\mu, \theta^2)$

$\dfrac{\sum_{i=1}^{k}(z_i - \bar{z})^2}{\theta^2} \sim \chi^2_{k-1}$

$\Rightarrow \sum_{i=1}^{k}(z_i - \bar{z})^2 \sim \theta^2 \chi^2_{k-1}$

$\eta + \tau_i + \bar{\varepsilon}_{i\cdot} \quad \eta + \bar{\tau}_\cdot + \bar{\varepsilon}_{\cdot\cdot}$

$$\bar{y}_{i\cdot} - \bar{y}_{\cdot\cdot} = \left(\tau_i - \bar{\tau}_\cdot\right) + \left(\bar{\varepsilon}_{i\cdot} - \bar{\varepsilon}_{\cdot\cdot}\right)$$

LNp.3-5

$$SSTr = \sum_{i=1}^{k} n\left(\bar{y}_{i\cdot} - \bar{y}_{\cdot\cdot}\right)^2$$

— independent —

$$= n\left\{ \sum_{i=1}^{k}\left(\tau_i - \bar{\tau}_\cdot\right)^2 + \sum_{i=1}^{k}\left(\bar{\varepsilon}_{i\cdot} - \bar{\varepsilon}_{\cdot\cdot}\right)^2 + 2\sum_{i=1}^{k}\left(\bar{\varepsilon}_{i\cdot} - \bar{\varepsilon}_{\cdot\cdot}\right)\left(\tau_i - \bar{\tau}_\cdot\right)\right\}.$$

mean = 0

The cross product term has mean 0 (because $\tau$ and $\varepsilon$ are independent). It can be shown that

$\sim \sigma_\tau^2 \chi^2_{k-1}$    $\sim \dfrac{\sigma^2}{n}\chi^2_{k-1}$

$$E\left(\sum_{i=1}^{k}\left(\tau_i - \bar{\tau}_\cdot\right)^2\right) = (k-1)\sigma_\tau^2 \quad \text{and} \quad E\left(\sum_{i=1}^{k}\left(\bar{\varepsilon}_{i\cdot} - \bar{\varepsilon}_{\cdot\cdot}\right)^2\right) = \frac{(k-1)\sigma^2}{n}.$$

iid                average of $\tau_i$'s          iid                average of $\bar{\varepsilon}_{i\cdot}$'s
$N(0, \sigma_\tau^2)$                                    $N(0, \sigma^2/n)$

Therefore

$$E(SSTr) = n(k-1)\sigma_\tau^2 + (k-1)\sigma^2, \qquad \text{cf.} \quad \text{In FEM (LNp. 3-6)}$$

$$E(MSTr) = E\left(\frac{SSTr}{k-1}\right) = \sigma^2 + n\sigma_\tau^2.$$

$E(SSTr) = n \cdot \sum_i (\tau_i - \bar{\tau})^2 + (k-1)\sigma^2$

---

# Variance components: estimation of $\sigma^2$ and $\sigma_\tau^2$

- From equations (1) and (2) in LNp. 3-35, we obtain the following unbiased estimates of the variance components:

Can this be always $\geq 0$?
(Note. $\sigma_\tau^2 \geq 0$)

same as the $\hat{\sigma}^2$ in FEM

$$\hat{\sigma}^2 = MSE \quad \text{and} \quad \hat{\sigma}_\tau^2 = \frac{MSTr - MSE}{n}.$$

check LNp.3-35    $MSTr/MSE$

Note that $\hat{\sigma}_\tau^2 \geq 0$ if and only if $MSTr \geq MSE$, which is equivalent to $F \geq 1$. Therefore a *negative* variance estimate $\hat{\sigma}_\tau^2$ occurs only if the value of the $F$ statistic is less than 1. Obviously the null hypothesis $H_0$ is not rejected when $F \leq 1$. Since variance cannot be negative, a negative variance estimate is replaced by 0. This does not mean that $\sigma_\tau^2$ is zero. It simply means that there is not enough information in the data to get a good estimate of $\sigma_\tau^2$.

$E(F_{n_1, n_2}) = \dfrac{n_2}{n_2 - 2}$

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

not "accept $H_0$"

- For the pulp experiment, $n = 5$, $\hat{\sigma}^2 = 0.106$, $\hat{\sigma}_\tau^2 = (0.447 - 0.106)/5 = 0.068$, i.e., sheet-to-sheet variance (within same operator) is 0.106, which is about 50% higher than operator-to-operator variance 0.068.

cf.

a property of operator population

*Implications on process improvement*: try to reduce the two sources of variation, also considering costs.

**Estimation of Overall Mean** $\eta$ ← the only fixed effect in REM

**check graph in LNp 3-31**

the intercept parameter in FEM is usually of **no** interest

- In REM, $\eta$, the population mean, is often of interest.

  From $E(y_{ij}) = \eta$, we use the estimate

  cf.

**In FEM, $E(y_{ij}) = \mu_i = \eta + \tau_i$**   cf.

**For balanced data, GLS = OLS**

$$\hat{\eta} = \bar{y}_{..}$$

Same as the $\hat{\eta}$ in FEM under sum coding, but in the case of FEM $\eta = (\mu_1 + \cdots + \mu_k)/k$

- $Var(\hat{\eta}) = Var(\bar{\tau}_{.} + \bar{\varepsilon}_{..}) = \frac{\sigma_\tau^2}{k} + \frac{\sigma^2}{N}$, where $N = \sum_{i=1}^k n_i$.

$\hat{\eta} = \bar{y}_{..} = \eta + \bar{\tau}_{.} + \bar{\varepsilon}_{..} \sim N(\eta, \sigma_\tau^2/k + \sigma^2/N)$       → $= E(MSTr)$

**pivotal quantity** $\dfrac{\bar{y}_{..} - \eta}{se(\hat{\eta})}$

For $n_i = n$, $Var(\hat{\eta}) = \frac{\sigma_\tau^2}{k} + \frac{\sigma^2}{nk} = \frac{1}{nk}\left(\sigma^2 + n\sigma_\tau^2\right)$.

$s.e.(\hat{\eta}) = \sqrt{\dfrac{MSTr}{nk}}$

Using (2) in LNp.3-35, $\frac{MSTr}{nk}$ is an unbiased estimate of $Var(\hat{\eta})$.

Confidence interval for $\eta$:   $\bar{y}_{..}$ and MSTr are indep. (LNp 33)

$SSTr \sim (\sigma^2 + n\sigma_\tau^2)\chi_{k-1}^2$

**estimate ± (critical value) × s.e. (estimate)**

$$\hat{\eta} \pm t_{k-1,\frac{\alpha}{2}}\sqrt{\frac{MSTr}{nk}}$$    cf.    ← N

**In FEM (under sum coding) C.I. for $\eta$:** $\bar{y}_{..} \pm t_{N-k,\frac{\alpha}{2}}\sqrt{\dfrac{MSE}{N}}$

- In the pulp experiment, $\hat{\eta} = 60.40$, $MSTr = 0.447$, and the 95% confidence interval for $\eta$ is

**compare REM and Split-plot design (LNp 4-45~66, future lecture)**

$SSE \sim \sigma^2\chi_{N-k}^2$

$$60.40 \pm 3.182\sqrt{\frac{0.447}{5 \times 4}} = [59.92, 60.88].$$

❖ **Reading**: textbook, 2.5