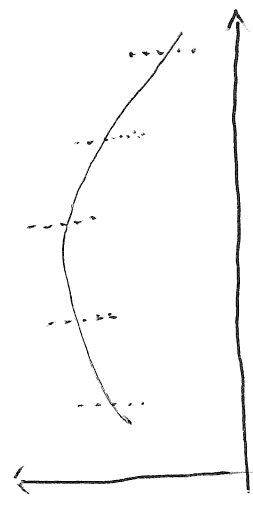Analysis of factor effects (52)

Then one can try to see if there is a "response function"
relating predictor $x$ and response $y$.

Ex: $x = $ price.    $y = $ sales.



(1). plot and see (dot plot)

try different curves

(2) Then test for lack of fit

$H_0: Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \varepsilon_i$

ANOVA

| Regression | | |
|---|---|---|
| SSR | SSTR | $r-1$ |
| | $p-1$ | |
| SSE | SSPE | $n_T - r$ |
| | $n_T - p$ | |
| SSTO | SSTO | $n_T - 1$ |
| | $n_T - 1$ | |

$SSLF = SSTR - SSR = SSE - SSPE$

$F_{LOF} = \dfrac{MSLF}{MSPE} = \dfrac{SSLF/(r-p)}{SSPE/(n_T-r)}$

Ex:
Page 771

# Chapter 18   ANOVA Diagnostics

(1). Examine if the model is appropriate ( check assumptions )

      indep. constant variance

      Normal

(2). If not, consider remedial measures.

(3) after that, repeat inference

---

① Analysis of residuals.

    Assumption :   $\varepsilon_{ij} \sim iid \; N(0, \sigma^2)$

        $e_{ij} := y_{ij} - \bar{y}_{i\cdot}$

- standardized residual    ( divided by an est. of $\sigma$ )

    $e_{ij}^* = \dfrac{e_{ij}}{\sqrt{MSE}}$

    $e_{ij}^*$ approximately $N(0,1)$.

    $|e_{ij}^*| > 3$ are potential outliers

• studentized residuals

$$\text{Var}(e_{ij}) = \text{Var}(y_{ij}) + \text{Var}(\bar{y}_{i\cdot}) - 2\,\text{Cov}(y_{ij}, \bar{y}_{i\cdot})$$

$$= \sigma^2 + \frac{\sigma^2}{n_i} - 2\frac{\sigma^2}{n_i} = \sigma^2(1 - \frac{1}{n_i})$$

$$S(e_{ij}) = S\sqrt{1 - \frac{1}{n_i}}$$

$$r_{ij} = \frac{e_{ij}}{S\sqrt{1 - \frac{1}{n_i}}}$$

• Studentized deleted residuals

(delete case $ij$, and refit the model,
compute the predicted value and residual
for the case $ij$)

$$t_{ij} = \frac{e_{ij}}{\sqrt{\dfrac{SSE(1 - \frac{1}{n_i}) - e_{ij}^2}{n_T - r - 1}}}$$

• both $\bar{y}_{i\cdot}$ and $S^2_{(-ij)}$ excludes $y_{ij}$

The new $\bar{y}_{i\cdot(-ij)} = \dfrac{n_i\,\bar{y}_{i\cdot} - y_{ij}}{n_i - 1}$, $\quad e_{ij'} = y_{ij'} - \bar{y}_{i\cdot(-ij)}$

· Outliers test

· We are selecting the further outlier, it is not legitimate to use a simple t-test.

· Bonferroni adjustment

Ex: if 5 observations, in 2 factor level(s)

By the Bonferroni method, compare with

$$t_{\frac{.05}{2 \times 5}} = t_{.005} = 9.925 \quad (\text{for } 5 - 2 - 1 = 2 \text{ df})$$

Tools of diagnotics

Residual plots

$r_{ij}$ vs $\bar{y}_{i.}$

$r_{ij}$ aligned dot plots , Normal prob. plot

They show · non constant variance $\underline{\vdots \cdots}$

· non independent errors (in a time plot)

· outliers (Bonferroni, check 2-3 points)

· nonnormality

Test for constant variance. (5-6)

$H_0: \sigma_1^2 = \sigma_2^2 = \cdots = \sigma_r^2$   vs   $H_A$ : $H_0$ not holed

① Hartley test.

$$H_{obs} = \frac{\max S_i^2}{\min S_i^2} ,$$   one-sided test

it has H statistic distr. under $H_0$ with parameters

$(r, n)$ :   $r$ : # of levels

   $n$ : Size of each group ( need same $n_i$ )

   $n_i \approx n \Rightarrow$ still ok.

Use $H_{(1-\alpha, r, n-1)}$ right-sided.

② Modified Levene test ( Brown - Forsythe test )

   $\widetilde{y}_{i\bullet} = i^{th}$ median = median of $(y_{i1}, \ldots, y_{in_i})$

   $d_{ij} = | y_{ij} - \widetilde{y}_{i\bullet} |$

Idea : Under $H_0$, $E(d_{ij}) = const$

   $\Rightarrow$ do ANOVA for $\{ d_{ij} \}$

$$MSTR = \frac{\Sigma\, n_i (\bar{d}_{i.} - \bar{d}_{..})^2}{r-1}$$

$$MSE = \frac{\Sigma\Sigma\, (d_{ij} - \bar{d}_{i.})^2}{n_T - r}$$

$$F_L^* = \frac{MSTR}{MSE} \approx F_{(r-1,\; n_T - r)}$$

Under $H_0$, $F_L^* \approx F_{(r-1,\, n_T - r)}$

Levene's — Use $\bar{y}_{i.}$ mean

R command leveneTest( )

Ex: 6 patients with 3 treatment, the responses were recorded

| A | 30 | 40 |
|---|----|----|
| B | 20 | 50 |
| C | 10 | 60 |

Use the modified Levene test to test a const. variance

| Treatment | $y_{ij}$ | $\bar{y}_i$ | $d_{ij}=\lvert y_{ij}-\bar{y}_{i.}\rvert$ | $\bar{d}_{i.}$ | $d_{ij}-\bar{d}_{i.}$ | $\bar{d}_{i.}-\bar{d}_{..}$ |
|---|---|---|---|---|---|---|
| A | 30 40 | 35 | 5 5 | 5 | 0 0 | -10 |
| B | 20 50 | 35 | 15 15 | 15 | 0 0 | 0 |
| C | 10 60 | 35 | 25 25 | 25 | 0 0 | 10 |

Test $H_0 : \sigma_1 = \sigma_2 = \sigma_3$ , $H_A :$ $H_0$ doesn't hold.

Carry out an ANOVA F based on $d_{ij}$

Get $SSTR = \sum_i n_i (\bar{d}_{i.} - \bar{d}_{..})^2 > 0$ , $r - 1 = 2$

$SSE = \sum \sum (d_{ij} - \bar{d}_{i.})^2 = 0$ $n_T - r = 3$

$\Rightarrow F = \dfrac{MSTR}{MSE} = \dfrac{\not{0}}{\not{0}}$ , $H_0$ is rejected at

any sig-level

Comment:

If $n_i \leq 2$ for all $i$ , then $d_{ij} - \bar{d}_{i.} = 0$

So $SSE = 0$ $\Rightarrow F = +\infty$ always.

We hesitate to apply it.

③ Bartlett's test

doesn't require $n_1 = n_2 = \cdots n_r$.

$$\chi^2 = \frac{\sum\limits_{i=1}^{r} \log \left( \dfrac{S_{pooled}^2}{S_i^2} \right)^{n_i - 1}}{1 + \dfrac{1}{3(r-1)} \left( \sum \dfrac{1}{n_i - 1} - \dfrac{1}{\sum(n_i - 1)} \right)} \approx \chi^2_{r-1}$$

one-sided

Rule of thumb

ANOVA F test can tolerate non-constant variance to some extent.

It is usually fine

$$\frac{\max \{ S_i \}}{\min \{ S_i \}} \leq 2 \text{ or even } 3, \text{ especially}$$

when $n_i$ are roughly equal.

Remedies

$\Bigg\{$ Normal
$\sigma_i^2 \neq$ constant $\Rightarrow$ weighted least squares

$\Bigg\{$ non-Normal
$\sigma_i^2 \neq$ constant $\Rightarrow$ transform $y_{ij}$, box-cox

Don't help
or other departures $\Rightarrow$ Nonparametric test
from assumption

① Weighted least squares

$y_{ij} = \mu_i + \varepsilon_{ij}$ .  $\varepsilon_{ij} \sim N(0, \sigma_i^2)$

estimate $\sigma_i^2$ by $S_i^2$ and use weights $w_i = \frac{1}{S_i^2}$

$SSTR(w) = \sum_i w_i \, n_i (\bar{y}_{i\cdot} - \bar{y}_{\cdot\cdot})^2 = \sum_i n_i \left( \frac{\bar{y}_{i\cdot} - \bar{y}_{\cdot\cdot}}{S_i} \right)^2$

$SSE(w) = \sum_i w_i \sum_j (y_{ij} - \bar{y}_{i\cdot})^2 = \sum_i \frac{1}{S_i^2} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot})^2$

$\qquad\qquad = \sum_i (n_i - 1) = n_T - r$

Then $F_w = \dfrac{MSTR(w)}{MSE(w)} \sim F_{(r-1, \, n_T - r)}$

② Transformations

variance - stabilizing

$\sigma^2 \sim \mu$ (poisson) $\Rightarrow y' = \sqrt{y}$ or $y' = \sqrt{y} + \sqrt{y+1}$

$\sigma^2 \sim \mu^2 \qquad\Rightarrow y' = \log y$

$\sigma^2 \sim \mu^4 \qquad\Rightarrow y' = \frac{1}{y}$

Generally $\sigma \sim \mu^\alpha$

$\lambda \longrightarrow \begin{cases} y^{1-\alpha} & \alpha \neq 1 \\ \log y & \alpha = 1 \end{cases}$

Box-cox

$$y' = y^\lambda, \quad \text{search for best } \lambda \text{ numerically}$$
to minimize SSE

Drawback:

• Except for a few special transformations ($\log, \sqrt{}, \frac{1}{y}$)
the transformed response lacks natural interpretation.

• Transformation doesn't work (helps little) for symmetric
but heavy-tailed distr. (many outliers)

• ANOVA F test is robust to non-normality, but it is
not resistant to outliers.

• If outliers are cannot be removed, try non-parametric
test.

③ Nonparametric Rank F - test.

Def : $R_{ij} = $ rank of $y_{ij}$ among all $n_T$ obs.

( $R_{ij} = r \iff y_{ij}$ is the $r$th smallest

If $y_{ij} = y_{i'j'} = y_{i''j''}$ etc $\Rightarrow$ average ranks)

Then do ANOVA of ranks

$$R_{ij} = \mu_i^{(R)} + \varepsilon_{ij}^{(R)}$$

$$F_R = \frac{MSTR^{(R)}}{MSE^{(R)}} \sim F_{r-1,\ n_T-r} \quad \text{if } n_i$$

are not very small.

Ex:

| Treatment | observations | | |
|---|---|---|---|
| 1 | -5 | 1 | 9 |
| 2 | -4 | 3 | 11 |
| 3 | 0 | 6 | 17 |

$\Rightarrow$

| | Ranks | | |
|---|---|---|---|
| 1 | 4 | 7 | 4 |
| 2 | 5 | 8 | 5 |
| 3 | 6 | 9 | 6 |

$$\bar{R}_{\cdot\cdot} = \frac{n_T+1}{2} = 5$$

$$SSTR = \sum_i n_i \left(\bar{R}_{i\cdot} - \bar{R}_{\cdot\cdot}\right)^2$$

$$= 3(4-5)^2 + 3(5-5)^2 + 3(6-5)^2 = 6 \quad \text{with dfs } 2$$

$$SSE = \sum_i \sum_j \left(R_{ij} - \bar{R}_{i\cdot}\right)^2$$

$$= (1-4)^2 + (4-4)^2 + (7-4)^2 + (2-5)^2 + (5-5)^2 + (8-5)^2 + (3-6)^2 + (6-6)^2 + (9-6)^2$$

$$= 54 \quad \text{with df } 6$$

$$F_R = \frac{SSTR/2}{SSE/6} = \frac{3}{9} = \frac{1}{3}$$

⑬

Kruskal - Wallis Test

$H_0: \mu_1 = \mu_2 = \cdots \mu_r$

Test statistic.

$$\chi^2_{kw} = \frac{SS_{TR} \ (\text{based on ranks})}{SS_{Total} \ (\text{based on ranks})/(n_T - 1)}$$

$$\approx \chi^2_{r-1} \quad \text{under } H_0 \quad \text{for large } n.$$

$$SS_{Total} = \sum_i \sum_j (R_{ij} - \bar{R}_{..})^2$$

$$= \sum_{k=1}^{n_T} (k - \frac{n_T + 1}{2})^2$$

$$= \frac{n_T (n_T^2 - 1)}{12}$$

$$\chi^2_{kw} = \frac{SS_{TR}(R)}{n_T (n_T + 1)/12}$$

Two way ANOVA

Consider two factors. A and B

Treatment = each combination of a level of A and a level of B

$$Ex \begin{cases} a = 2 & \text{levels of A} \\ b = 3 & \text{levels of B} \end{cases} 6 \text{ treatment}$$

Say we have $n_T = 36$ obs. 6 for each treatment.

This is a complete design.

If not all, but only a fraction of treatments is used in study $\Rightarrow$ fractional design.

2- way ANOVA model.

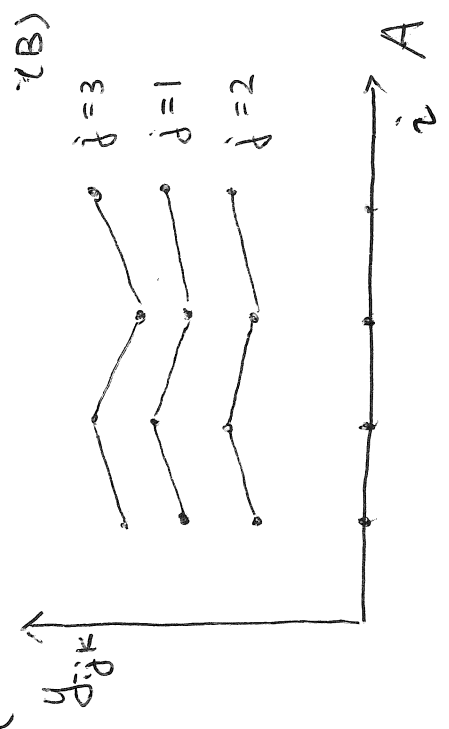cell mean $\mu_{ij} = E(Y)$ ( ith level of A )
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ jth level of B

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad i = 1, \dots a$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad j = 1, \dots b$

$$\mu_{i\cdot} = \frac{1}{b}\sum_{j=1}^{b}\mu_{ij}$$

$$\mu_{\cdot j} = \frac{1}{a}\sum_{i=1}^{a}\mu_{ij}$$

$$\mu_{\cdot\cdot} = \frac{1}{ab}\sum_{i}\sum_{j}\mu_{ij}$$

main effects

$$\alpha_i = \mu_{i\cdot} - \mu_{\cdot\cdot} \Rightarrow \sum_{i=1}^{a}\alpha_i = 0$$

$$\beta_j = \mu_{\cdot j} - \mu_{\cdot\cdot} \Rightarrow \sum_{j}\beta_j = 0$$

Additive factor effects

$$\Leftrightarrow \mu_{ij} = \mu_{\cdot\cdot} + \alpha_i + \beta_j$$

$$= \mu_{\cdot\cdot} + \mu_{i\cdot} - \mu_{\cdot\cdot} + \mu_{\cdot j} - \mu_{\cdot\cdot} = \mu_{i\cdot} + \mu_{\cdot j} - \mu_{\cdot\cdot}$$

Then for any $i, j, k, \ell$

$$\mu_{ij} + \mu_{k\ell} = 2\mu_{\cdot\cdot} + \alpha_i + \beta_j + \alpha_k + \beta_\ell$$

$$\mu_{i\ell} + \mu_{kj} = 2\mu_{\cdot\cdot} + \alpha_i + \beta_\ell + \alpha_k + \beta_j$$

$$\Rightarrow \mu_{ij} - \mu_{i\ell} = \mu_{kj} - \mu_{k\ell} \qquad \text{no interaction}$$

( Doesn't matter if (i) paired with (j) or (l) )
effect of A doesn't depend on the level of B )

(16)

Illustration

(AB)

$y_{ijk}$

$j=3$

$j=1$

$j=2$

$i$ A

$\longleftarrow$ treatment means plot

Non-additive (interacting) factor effects

$y_{ijk}$

$i=1$
$i=2$

interaction effect is

$$(\alpha\beta)_{ij} := \mu_{ij} - (\mu_{..} + \alpha_i + \beta_j)$$

$$= \mu_{ij} - \alpha_i - \beta_j - \mu_{..}$$

$$= \mu_{ij} - \mu_{i.} - \mu_{.j} + \mu_{..}$$

∃ interaction if

(1) $\mu_{ij} \neq \mu_{..} + \alpha_i + \beta_j$ for some $i, j$

(2) $\mu_{ij} - \mu_{il} \neq \mu_{kj} - \mu_{kl}$ for some $i, j, k, l$

(3) treatment means curves are not parallel.

Some $(\alpha\beta)_{ij} = 0$ is possible

always $\sum_i (\alpha\beta)_{ij} = \sum_j (\alpha\beta)_{ij} = 0$

$$\sum_i (\alpha\beta)_{ij} = \sum_i (\mu_{ij} - \alpha_i - \beta_j - \mu_{..})$$

$$= a\,\mu_{.j} - 0 - a\beta_j - a\mu_{..}$$

· Transformable interactions.

Say $\mu_{ij} = \mu_{..}\, \alpha_i\, \beta_j$

$\log \mu_{ij} = \log \mu_{..} + \log \alpha_i + \log \beta_j$

or $\mu_{ij} = \alpha_i + \beta_j + 2\sqrt{\alpha_i \beta_j}$

$\sqrt{\mu_{ij}} = \sqrt{\alpha_i} + \sqrt{\beta_j}$