

Analysis of Messy Design

“Does it Matter?”

Messy Data



©1999 HOLD THE MUSTARD PRODUCTIONS

Agenda

- “Messy Data”
- Why do Analysts Care?
- Example
 - Single Treatment ANOVA w/ Heterogeneous Variance Terms
- Learning Points
- References: Analysis of Messy Data: Volume I, Design of Experiments, George A. Milliken and Dallas E. Johnson. Chapman and Hall, New York, 1992.

What is Messy Data?

- Failure of statistical assumptions that influence underlying distribution theory:
 - Homogeneity of variance across treatments
 - Severe outliers in data sets
 - Missing treatments and/or treatment combinations in an experiment
 - Unequal number of observations in experiments

Why Do Analysts Care?

- Experiments influence action: decisions on the influence of different treatments on a measured response
- Sometimes, we don't.
- Appropriate to study the differences.

The Model of Interest

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

Where, $\varepsilon_{ij} \square N(0, \sigma_i^2)$

The variance terms may be different for each population or treatment applied.

Try to address three questions:

- Are the treatment means equal?
- Can the analyst compare the different means using a pair-wise approach?
- Is the analyst able to use linear contrasts to test hypotheses of interest?

Distribution Results

For large n_i , let $Z = \frac{\left(\frac{\sum_i c_i \hat{\mu}_i - \sum_i c_i \mu_i}{\sqrt{\frac{\sum_i c_i \hat{\sigma}_i^2}{n_i}}} \right)}{\sqrt{\frac{\sum_i c_i \sigma_i^2}{n_i}}}$

The above statistic has an approximate $N(0,1)$ distribution.

Satterthwaite's Approximation

- If sample sizes are not large, then:

$$\hat{\nu} = \frac{\left(\sum_i c_i^2 \hat{\sigma}_i^4 / n_i \right)^2}{\sum_i \left[c_i^4 \sigma_i^4 / n_i^2 (n_i - 1) \right]}$$

Reject $H_0 : \sum c_i \mu_i = a$ with $|t_c| = \frac{\left| \sum_i c_i \hat{\mu}_i - a \right|}{\sqrt{\sum_i c_i \hat{\sigma}_i^2 / n_i}} > t_{\alpha/2, \hat{\nu}}$

Learning Points

- The data will guide analysis approach; assumptions are key
- Software applications lend themselves to analysis of the treatment model: SAS is a good starting point with GLM; discourage Minitab.