

# Designs with One Source of Variation

Ex) only one treatment factor

Hao Zhang

# Completely Randomized Design

- Used when we want to compare several treatments and the experimental units are homogeneous.
- No blocking or stratification.

Subjects

★ Definition

The experimental units are assigned to the treatments completely at random, subject to the number of replicates for each treatment.

The number of replicates is the number of units assigned to the treatment and equals the number of observations to be taken for the treatment.

★ Definition

# How to randomize

Ex) 3 treatments: A(1), B(2), C(3)

Replicates: 4 per treatment  $\rightarrow 4 \times 3 = 12$  total # of obs.

Notation:  $r_i$  to denote the number of replicates for  $i$ th treatment, and  $n = \sum r_i$ . Code the treatments from 1 to  $v$  and label the experimental units 1 to  $n$ .  
 $v$ : # of treatments  
 $n$ : total # of obs.

- ▶ Step 1: Enter into one column  $r_1$  1's, then  $r_2$  2's, ..., and finally  $r_v$   $v$ 's, giving a total of  $n = \sum r_i$  entries. These represent the treatment labels. 1,1,1,1,2,2,2,2,3,3,3,3  
repeated treatment labels
- ▶ Step 2: Enter into another column  $n = \sum r_i$  random numbers, including enough digits to avoid ties. 12 random #'s
- ▶ Step 3: Reorder both columns so that the random numbers are put in ascending order. This arranges the treatment labels into a random order.
- ▶ Step 4: Assign experimental unit  $t$  to the treatment whose label is in row  $t$ .

This ensures treatments are randomly distributed to units.

# Model for CRD

Let the experiment have  $v$  treatments and the  $i$ th treatment has  $r_i$  replicates. Let  $Y_{it}$  denote the response obtained on the  $t$ th observation of the  $i$ th treatment.

$$Y_{it} = \mu_i + \epsilon_{it}, \quad t = 1, \dots, r_i, \quad i = 1, \dots, v$$

$\mu_i = \text{Mean effect of treatment } i$

where  $\epsilon_{it} \sim N(0, \sigma^2)$  are all independent (i.i.d.)  
 $\downarrow$  random noise

Model parameters:  $\mu_i, i = 1, 2, \dots, v$  and  $\sigma^2$ .

Equivalent treatment effects model:

$$Y_{it} = \mu + \tau_i + \epsilon_{it}, \quad t = 1, \dots, r_i, \quad i = 1, \dots, v$$

$\tau_i = \text{treatment effect}$

There are numerous ways to write  $\mu_i = \mu + \tau_i$ . Hence,  $\tau_i$ 's are not uniquely specified. Unless we impose constraint  $\downarrow$

- ▶  $\mu = (1/v) \sum_i^v \mu_i, \quad \tau_i = \mu_i - \mu.$
- ▶  $\tau_v = 0, \quad \mu = \mu_v, \quad \tau_i = \mu_i - \mu_v.$  SAS implements this one.

# Estimability

Because  $\tau_i$ 's are not unique, they are not estimable. However, if  $\sum_i c_i = 0$ ,  $\sum c_i \tau_i = \sum c_i \mu_i$  as long as  $\mu_i = \mu + \tau_i$  for all  $i$ . It is usually called a contrast that is estimable. Any linear combination  $\sum c_i \mu_i$  such that  $\sum c_i = 0$ .

## Examples:

- ▶  $\tau_1 - \tau_2$  "differences between treatment 1 and 2."
- ▶  $(\tau_1 + \tau_2)/2 - \tau_3$  "is the average of 1 and 2 better than 3?"
- ▶  $(\tau_1 + \tau_2)/2 - (\tau_3 + \tau_4)/2$  "Comparing two groups of treatments."

Everything expressed in terms of treatment means are estimable!

# Parameter Estimation

# Notations

- ▶  $Y_{it}$ : the random variable representing the  $t$ th outcome of treatment  $i$ .
- ▶  $y_{it}$ : the observed value of  $Y_{it}$ .
- ▶ I might suppress the differences in my notes.



Sample mean for treatment  $i$

$$\bar{y}_{i.} = (1/r_i) \sum_{t=1}^{r_i} y_{it}, \bar{y}_{..} = \frac{1}{n} \sum_{i=1}^v \sum_{t=1}^{r_i} y_{it} = \frac{1}{n} \sum_{i=1}^v y_{i.} = \frac{1}{n} y_{..} = \bar{y}_{..}$$

where  $n = \sum_{i=1}^v r_i$ .

Overall mean

# Least Squares Estimation: (LSE)

Objective:

Minimize the sum of squared errors (SSE)

$$SSE = \sum_{i=1}^v \sum_{t=1}^{r_i} (y_{it} - \mu_i)^2.$$

Solution:

$$\hat{\mu}_i = \bar{y}_{i..}$$

★ Sample mean of each treatment is the best estimate of its true mean.



# Properties:

- ▶ The LSE is the best linear unbiased estimator (BLUE).  
 $E(\bar{Y}_{i.}) = \mu_i$  (unbiased).
- ▶  $\bar{Y}_{i.} \sim N(\mu_i, \sigma^2/r_i)$ . *Normally distributed*
- ▶ For any constants

$$\sum_{i=1}^v c_i \bar{Y}_{i.} \sim N\left(\sum_{i=1}^v c_i \mu_i, \sigma^2 \sum_{i=1}^v \frac{c_i^2}{r_i}\right)$$

*supports hypothesis testing + confidence intervals.*

## Estimation of $\sigma^2$ : Residual Variance

$$SSE = \sum_i \sum_t (Y_{it} - \bar{Y}_{i.})^2 = \sum_i (r_i - 1) \underbrace{S_i^2}_{\text{Sample variance}}$$

$\hat{\sigma}^2 = SSE / (n - v) = \underbrace{MSE}_{\text{unbiased estimator of } \sigma^2} = \text{Variance within treatments}$

$$E(\hat{\sigma}^2) = \sigma^2.$$

**Upper Confidence Limit of  $\sigma^2$**

# Upper Confidence Limit of $\sigma^2$

$SSE/\sigma^2$  has a  $\chi^2$ -distribution with  $(n - v)$  degrees of freedom.  
Hence, the 95% upper confidence limit for  $\sigma^2$

$$\frac{SSE}{\chi_{n-v,0.95}^2}$$

where  $\chi_{n-v,0.95}^2$  is the 5<sup>th</sup> percentile of the chi-square distribution with  $(n - v)$  degrees of freedom.

**SAS Code**

# Percentiles of distribution

The SAS program to calculate the 95% percentile of a chi-square distribution with 13 degrees of freedom is given below.

```
data chisq;  
input prob df;  
  percentile=cinv(prob, df);  
  lines;  
  0.05 13  
;  
proc print data=chisq;  
run;
```

# The Prius Experiment

In an experiment to study the effects of drivers on the mpg of Toyota Prius, 12 new Prius were randomly assigned to three drivers so that each driver drove four cars and obtained the mpgs. This is a completely randomized design. The data are given below.

<i>d1</i>	<i>d2</i>	<i>d3</i>
50.33	48.11	49.08
46.83	50.14	48.89
51.57	43.22	49.96
45.33	47.26	49.70

The SAS program to get the sample means and sample standard deviations:

# SAS Code

```
data prius;  
input driver mpg;  
lines;  
1 50.33  
1 46.83  
1 51.57  
1 45.33  
2 48.11  
2 50.14  
2 43.22  
2 47.26  
3 49.08  
3 48.89  
3 49.96  
3 49.70  
;  
run;  
proc print data=prius;  
run;
```



# SAS Code

```
proc means data=prius;  
by driver;  
run;
```

# Estimation

1. Find an estimate of the variance  $\sigma^2$ . The estimate is given by the MSE:

$$MSE = \frac{1}{n - v} \sum_{i=1}^3 (r_i - 1) s_i^2 \quad (1)$$

$$= \frac{1}{\underset{n}{12} - 3} (3 \underset{s_1^2}{*} 2.92^2 + 3 \underset{s_2^2}{*} 2.90^2 + 3 \underset{s_3^2}{*} 0.51^2) \quad (2)$$

$$= 5.75. \quad (3)$$

2. Find the 95% confidence upper limit for  $\sigma^2$ . First we need to find the 5<sup>th</sup> percentile for the  $\chi^2$ -distribution with  $\overset{12}{n} - \overset{3}{v} = 9$  degrees of freedom, which equals  $\chi_{9,0.95}^2 = 3.325$ . The 95% confidence upper limit is given by  $\frac{SSE}{\chi_{9,0.95}^2} = \frac{\overset{df}{9} * \overset{MSE}{5.75}}{3.325} = 15.55$ .