



# 南京师范大学附属中学

HIGH SCHOOL AFFILIATED TO NANJING NORMAL UNIVERSITY

- $\binom{n}{k} = \frac{n!}{k!(n-k)!}$
- randomized experiment
- assumption: ① no difference between 2 treatments ② ...
- randomization p-value  
p-value: prob of getting a test statistic as extreme or more extreme than the observed value of test statistics  $t^*$   
 $P(T \leq t^* | H_0) = \sum \frac{I(t_i \leq t^*)}{\binom{N}{N_A}}$  (one sided)
- two sided (not simply double prob of one tail b/c of non-symmetry)  
 $P(|T - E| \geq |t^* - E| | H_0) = \sum \frac{I(|t_i - E| \geq |t^* - E|)}{\binom{N}{N_A}}$
- use median to replace mean works as well (also for one/two sided)
- why? mean is affected by extreme data value easily, but median better reflects the "middle class" data feature.
- Monte Carlo sampling (when  $\binom{30}{15} \rightarrow$  large), practical unless size is small)
  - randomization
  - one side p-value:  $p = \frac{1 + \sum_{i=1}^M I(t_i \geq t^*)}{M+1}$  including the observed  $t^*$  there are  $M+1$  test statistics.
- decision theory
 

$H_0$ true	$H_1$ true	
Accept $H_0$	correct decision	type II error
Reject $H_0$	type I error	correct decision

$P(\text{type I}) = \alpha$  power:  $1 - \beta$  rejecting  $H_0$  when  $H_1$  is true.
- exact or conservative, test is guaranteed to control the prob of type I error under very minimal conditions: randomizations of experimental units to the treatments.
- two-sample t-test if  $\mu_A, \mu_B \sim N(\mu, \sigma^2)$  with same  $\sigma^2$
- ~~$\bar{Y}_A - \bar{Y}_B \sim N(\mu_A - \mu_B, \sigma^2(\frac{1}{n_A} + \frac{1}{n_B}))$ ,  $\frac{\bar{Y}_A - \bar{Y}_B - \delta}{\sigma \sqrt{\frac{1}{n_A} + \frac{1}{n_B}}} \sim N(0, 1)$ ,  $\delta = \mu_A - \mu_B$~~
- $$\frac{\bar{Y}_A - \bar{Y}_B - \delta}{S \sqrt{\frac{1}{n_A} + \frac{1}{n_B}}} \sim t_{n_A + n_B - 2}$$
 two-sample t-test       $t^* = \frac{\bar{Y}_A - \bar{Y}_B}{S \sqrt{\frac{1}{n_A} + \frac{1}{n_B}}}$
- ~~What about qq plot? if data points along qq line  $\Rightarrow$  normality is satisfied!~~
- randomization & 2-sample t-test have almost identical p-values.  
 $\Rightarrow$  exchangeability is true!
- randomized paired design. (flip a coin!) comparison

$$P(D \leq d^* | H_0) = \sum_{i=1}^N \frac{I(d_i \leq d^*)}{N}$$

$d^*$ : observed diff

D: A - B



# 南京师范大学附属中学

HIGH SCHOOL AFFILIATED TO NANJING NORMAL UNIVERSITY

- When # of ways to split data is large, use MC.
- Experiment (if assignment mechanism is controlled by researchers)
- If an subject receives both treatment in a paired fashion  $\Rightarrow$  r.p. design.

For clinical trial: we want to find if there is causation.

- randomization: ensure that the groups will be similar w.r.t. all factors measured in the study & all factors that are not measured.

Power and Sample size.

- one sample z-test: reject at  $\alpha$  iff  $|\frac{\bar{X}-\mu_0}{\sigma/\sqrt{n}}| \geq Z_{\alpha/2}$  (two-sided)

$$1-\beta = P(\text{Reject } H_0 | \mu = \mu_1) = P(|\bar{X} - \mu_0| \geq \frac{\sigma}{\sqrt{n}} Z_{\alpha/2} | \mu = \mu_1) = P(\bar{X} - \mu_0 \geq \frac{\sigma}{\sqrt{n}} Z_{\alpha/2} | \mu = \mu_1) + P(\bar{X} - \mu_0 \leq -\frac{\sigma}{\sqrt{n}} Z_{\alpha/2} | \mu = \mu_1)$$

$$= 1 - \Phi\left(Z_{\alpha/2} - \left(\frac{\mu_1 - \mu_0}{\sigma/\sqrt{n}}\right)\right) + \Phi\left(-Z_{\alpha/2} - \left(\frac{\mu_1 - \mu_0}{\sigma/\sqrt{n}}\right)\right),$$

~~prob~~ power: prob of z-test will detect a mean of  $\mu_1$ , when  $n=$ ,  $\sigma=$ ,  $\alpha=$  ..

- one sample t-test: reject at  $\alpha$  iff  $|\frac{\bar{X}-\mu_0}{S/\sqrt{n}}| \geq t_{n-1, \alpha/2}$ , similarly for  $1-\beta =$  ..

- two sample t-test:  $\bar{Y}_k = \frac{1}{n_k} \sum Y_{ik}$  sample mean.  $S_p^2 = \frac{1}{n_1+n_2-2} \sum_{k=1}^2 \sum_{j=1}^{n_k} (Y_{ik} - \bar{Y}_k)^2$

$$\text{2-sample t-statistic } T_n = \frac{\bar{Y}_1 - \bar{Y}_2}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1+n_2-2}$$

$$T_n \sim \underbrace{t_{n_1+n_2-2} \text{ under } H_0}_{df} \text{ & } t_{n_1+n_2-2, \gamma} \text{ with noncentrality parameter } \gamma = \frac{\theta}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \text{ under } H_1$$

$$1-\beta = P(t_{n_1+n_2-2, \gamma} \geq t_{n_1+n_2-2, \alpha/2}) + P(t_{n_1+n_2-2, \gamma} \leq -t_{n_1+n_2-2, \alpha/2})$$

- power as a function of sample size & sd.

$$ES: \frac{\mu_1 - \mu_2}{\sigma}, SML: 0.2, 0.5, 0.8$$

type of allocations, { randomized  
nonrandomized

- sample size (known variance & equal allocation)

two-sample z-test --- (similar to two-sample t-test)

$$\text{solve } Z_\beta + Z_{\alpha/2} = \left( \frac{1\theta_1}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \right) \text{ assume } n_1 = n_2 = \frac{n}{2} \Rightarrow n = \frac{4\theta^2 (Z_\beta + Z_{\alpha/2})^2}{\theta^2}$$

$$\text{unequal allocation: } r = n_1/n_2 \Rightarrow n_1 = r \cdot n_2 \Rightarrow n_2 = \frac{(1 + \frac{1}{r}) \theta^2 (Z_\beta + Z_{\alpha/2})^2}{\theta^2}$$

imbalance  $\Rightarrow$  loss of power.

- Binary with proportion:  $n_2 = \frac{(Z_{\alpha/2} + Z_\beta)^2}{\theta^2} (P_1(1-P_1) \frac{1}{r} + P_2(1-P_2))$

- power by simulation:

- two sample test assumption: one is iid  $N(\mu_1, \sigma^2)$ , one is iid  $N(\mu_2, \sigma^2)$

- increase power? either increase  $n, \alpha$  or  $\theta$  or decrease  $\sigma$ .

- To increase power? either increase  $n, \alpha$  or  $\theta$  or decrease  $\sigma$ .

recommend  $n \uparrow$  b/c  $\alpha, \theta$  are set (stat.-meaningful)

$\sigma \downarrow$  is not practical.



# 南京師範大學附屬中學

HIGH SCHOOL AFFILIATED TO NANJING NORMAL UNIVERSITY

- causal inference
  - could have - zero causal effect but + predictive comparison
  - positive causal effect but 0 predictive comparison
- fundamental problem of causal inference.
  1. The definition of the causal ~~inference~~ effect depends on the potential outcomes, but not depends on which outcome is actually observed.
  2. The causal effect ~~of~~ is the comparison of potential outcome's, for the same unit at the same moment in time post-treatment.
- ∴ The fundamental problem of CI is at most one of 2 potential outcomes can be observed! Can NEVER measure a causal effect directly.
- Then HOW? ① close substitutes / ② Randomization & experiment / ③ statistical adjustment
- (2) do average treatment effect  $\bar{Y}' - \bar{Y}^0$
- Stable Unit Treatment Value Assumption (SUTVA): The potential outcomes for any unit do not vary with treatments assigned to other units, & for each unit, there are no different forms or versions of each treatment level, which lead to different outcomes.
- observational studies (not assigned, e.g. lung cancer vs. smoking), sometimes inside an experiment
- If treatment assignment  $T$  is conditionally independent of  $y^0, y^1$  given confounding covariates  $X \Rightarrow$  treatment assignment is ignorable.  $y^0, y^1 \perp T | X$
- ~~the "teaching"~~  
non-ignorable: e.g. teaching experience, teacher's motivation ...  
• If treatment assignments depend on info not included in the model then watch out!
- propensity score.  $e(\vec{x}) = P(T=1 | \vec{x})$ ,  $\vec{x}$  are obs. covariates. (before treatment)
  - ↳ known in experiments
  - ↳ unknown in obs. studies  $\Rightarrow$  estimated by models e.g. logistic regression.
$$\log \left( \frac{P_i}{1-P_i} \right) = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \hat{\beta}_3 x_{i3} \text{ where } p_i = P(T_i=1)$$

$$\hat{p}_i = \frac{\exp(\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots)}{1 + \exp(\dots)}$$
- balancing property of propensity score:  $P(\vec{x} | T=1, e(\vec{x})) = P(\vec{x} | T=0, e(\vec{x}))$  or  $T \perp \vec{x} | e(\vec{x})$
- ignorable treatment assignment  $P(T | Y(0), Y(1), \vec{x}) = P(T | \vec{x})$



# 南京師範大學附屬中學

HIGH SCHOOL AFFILIATED TO NANJING NORMAL UNIVERSITY

## ANOVA (analysis of variances)

- compare multiple samples.
- between treatment variation & within treatment variation

$$Y_{ij} - \bar{y}_{..} = (Y_{ij} - \bar{y}_{i.}) + (\bar{y}_{i.} - \bar{y}_{..}) \quad Y_{ij} \dots j\text{th obs under treatment } i=1, \dots, a$$

$$E(Y_{ij}) = \mu_i = \mu + T_i, \text{Var}(Y_{ij}) = \sigma^2. \quad T_i \text{ is } i\text{th treatment effect.}$$

$$\bar{y}_{i.} = \frac{\sum_j^n Y_{ij}}{n}, \quad \bar{y}_i = \bar{y}_{i.}/n$$

$$\bar{y}_{..} = \frac{\sum_i^a \sum_j^n Y_{ij}}{N}, \quad \bar{y}_{..} = \bar{y}_{..}/N$$

$$N = a \cdot n$$

$$H_0: \mu_1 = \dots = \mu_a \quad (\text{treatment})$$

$$H_a: \mu_i \neq \mu_j, \forall i \neq j$$

- ANOVA identity  
total sum of squares:  $SS_T = \sum_{i=1}^a \sum_{j=1}^n (Y_{ij} - \bar{y}_{..})^2 = n \sum_{i=1}^a (\bar{y}_{i.} - \bar{y}_{..})^2 + \sum_{i=1}^a \sum_{j=1}^n (Y_{ij} - \bar{y}_{i.})^2$

$\downarrow$   
 $df: N-1$

~~$SS_E = \sum_{i=1}^a \sum_{j=1}^n (Y_{ij} - \bar{y}_{i.})^2$~~

sample ~~i~~ith treatment  $S_i^2 = \frac{\sum_{j=1}^n (Y_{ij} - \bar{y}_{i.})^2}{n-1}$

Variance sum up all  $S_i^2 \rightarrow$  population var  $\frac{SS_E}{N-a}$  so  $SS_E$  is pooled est. of common var  $\sigma^2$  within mean square for treatment  $MS_{\text{treat}} = \frac{SS_{\text{treat}}}{a-1}$ , mean square for error:  $MS_E = \frac{SS_E}{N-a}$  each  $a$  treatments

$\frac{SS_{\text{treat}}}{\sigma^2} \sim \chi_{a-1}^2, \frac{SS_E}{\sigma^2} \sim \chi_{N-a}^2$ , thus for  $H_0: F = \frac{MS_{\text{treat}}}{MS_E} \sim F_{a-1, N-a}$ . (\*)

$$df: a-1 \quad \downarrow \quad df: a(n-1) = an-a = N-a$$

within each treatment

$n-1$  df, times  $a$   
 $= a(n-1)$

## General ANOVA

Source of variation	Df	SS	MS	F
B/w treatments	a-1	SS <sub>treat</sub>	MS <sub>treat</sub>	
Within treatments	N-a	SS <sub>E</sub>	MS <sub>E</sub>	$F = \frac{MS_{\text{treat}}}{MS_E}$

- Assumptions: ① additive model  $Y_{ij} = \mu + T_i + \varepsilon_{ij}$  (additive model)  
(behind p-values in F-test)  
↳ treatment effect  $T_i = \mu_i - \mu$

- ② errors  $\varepsilon_{ij}$  iid with common var  $\text{Var}(\varepsilon_{ij}) = \sigma^2 \forall i, j$ . (constant var)

$$E(MS_{\text{treat}}) = \sum_{i=1}^a T_i^2 + \sigma^2, E(MS_E) = \sigma^2$$

If no difference between treats means  $\Rightarrow T_1 = \dots = T_4$ , both  $MS_E, MS_T$  estimate  $\sigma^2$

- ③ if  $\varepsilon_{ij} \sim N(0, \sigma^2)$  then  $MS_T$  &  $MS_E$  independent (normality)

Under  $H_0$  that  $\sum_{i=1}^a T_i^2 = 0$ ,  $F = \frac{MS_T}{MS_E}$  is ratio of 2 indepdnt est of  $\sigma^2 \Rightarrow$  (\*)

- Least square est.  $\hat{Y}_{ij} = \mu + T_1 X_1 + T_2 X_2 + T_3 X_3 + \varepsilon_{ij}$ ,  $X_i = \begin{cases} 1 & \text{if diet } B \\ 0 & \text{o.w.} \end{cases} \dots C, D$ .

# so lse:  $\hat{\mu} = \bar{Y}_{1.}, \hat{T}_1 = \bar{Y}_{2.} - \bar{Y}_{1.}, \hat{T}_2 = \bar{Y}_{3.} - \bar{Y}_{1.}, \hat{T}_3 = \bar{Y}_{4.} - \bar{Y}_{1.}$

— But NOT ASSUMPTIONS behind calculations for ss, df or MS.

# PRACFINAL

Date \_\_\_\_\_

2.

$$(a) \frac{1}{2}$$

$$(b) \lambda_1 f_1(x_1, x_2)$$

$$\lambda_2 f_2(x_1, x_2) = \frac{1}{3} \cdot \frac{1}{\pi} + \frac{2}{3} \cdot \frac{2(x_1^2 + x_2^2)}{\pi} = \frac{1}{1 + 4(x_1^2 + x_2^2)}$$

$$0 < x_1^2 + x_2^2 < 1 \Rightarrow 1 > P > \frac{1}{5}$$

$$\frac{1}{5} < P < 1$$

Prac final.

Q1.

$$(a) \hat{R} = LL^T + \Psi$$

$$= \begin{pmatrix} 0.628 \\ 0.786 \\ 0.781 \\ 0.457 \\ 0.954 \end{pmatrix} (0.628 \cdots \rightarrow) + \begin{pmatrix} 0.605 & 0.383 & 0.389 & 0.084 & 0.089 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$= (5 \times 5)$$

appropriate since p-value = 0.761 > 0.05

$$(b) QQ^T = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = I$$

orthogonal.

$$L \rightarrow LQ$$

$$L = \begin{pmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{pmatrix} \quad LQ = ? \quad \checkmark$$

Q2.  $P(\text{missclassification}) = \frac{P(f_1, x_1 > 0)}{P(x_1 > 0)} = \frac{\frac{1}{3} \cdot \frac{1}{2}}{\frac{1}{2}} = \frac{1}{3}$

~~$P(\text{choose } x \text{ from } f_2 \mid x \text{ from } f_1, x_1 > 0)$~~

~~$P(\text{believe it from } f_2 \mid \text{actually from } f_1, x_1 > 0)$~~

~~$P(x_1 > 0 \mid \text{actually from } f_1, x_1 > 0)$~~

(b) Classify  $x$  to group  $j$  if  $\forall i \neq j$ .

$$\lambda_i f_i(x) < \lambda_j f_j(x)$$

$$\frac{f_j(x)}{f_i(x)} > \frac{\lambda_i}{\lambda_j}$$

$$\frac{f_2(x)}{f_1(x)} = \frac{2(x_1^2 + x_2^2)}{x_1^2} > \frac{\frac{1}{2}}{\frac{1}{3}} = \frac{\frac{1}{2}}{\frac{2}{3}} = \frac{3}{4} > \frac{1}{2}$$

$$\Rightarrow x_1^2 + x_2^2 > \frac{1}{4} \Rightarrow \text{Group 2.}$$

or, Group 1.

(c). Nothing b/c ~~no matter what signs~~  $x_1$  &  $x_2$  are, we are only interested in  $x_1^2$  &  $x_2^2$ , which must be ~~positive~~. Non-negative.  
So it's useless here.

### Q3. PCA vs ICA

~~same goal different results.~~

PCA  $\downarrow$  covariance of the data, ICA  $\downarrow$  higher order stats ~~etc~~  
e.g. Kurtosis, thus ~~more~~  $\downarrow$  the mutual info of the output,  
(what's left is "independent").

PCA outputs orthogonal vectors ~~resp~~ which are responsible for large portion of variances of the data. ICA identifies indepdnt components for non-normal data.

Besides, ICA has 2 ambiguities:

so ICA model equation is underdetermined system,

one cannot determine variances of the indpt components.

Also. cannot rank the order of dominant component.

Q4. How PC~~is~~ can be used to construct est. for factor analysis.

Q5.

(a). 2 & 3 in 1st clustering ✓

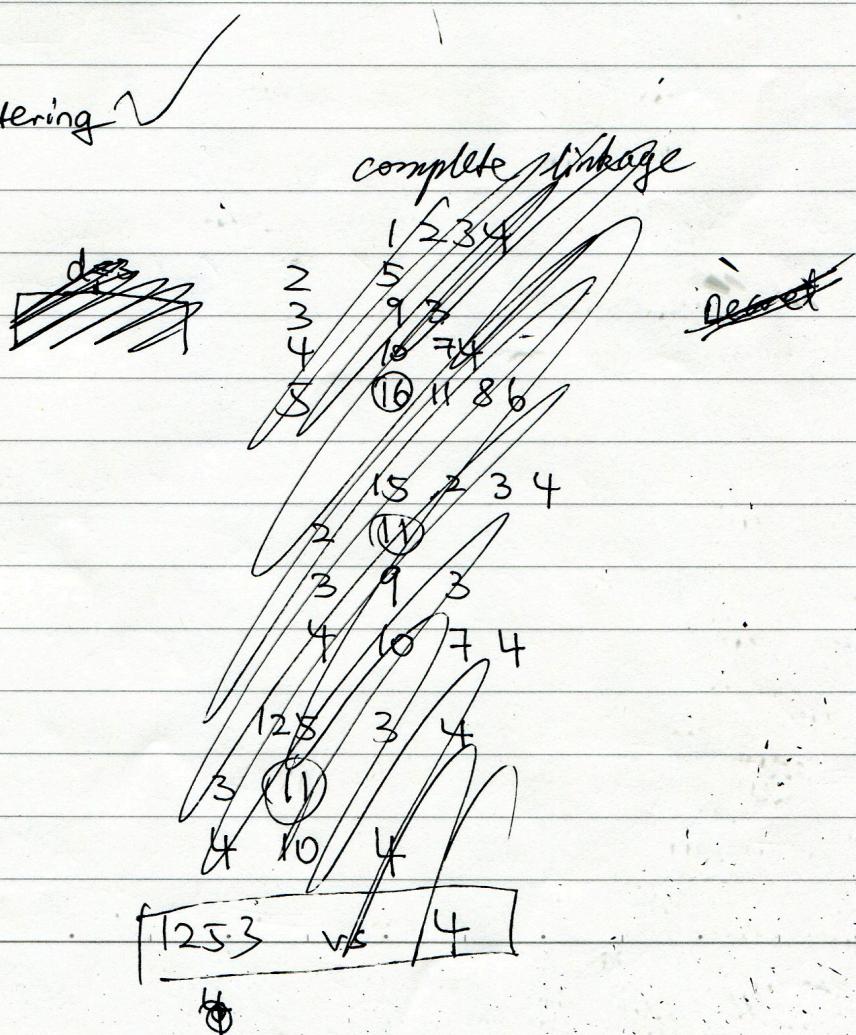
(b). Single linkage

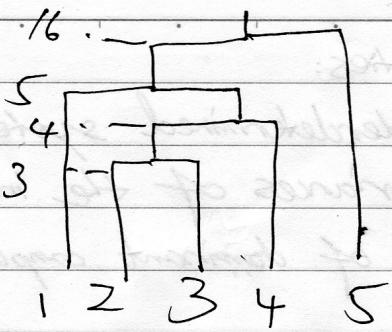
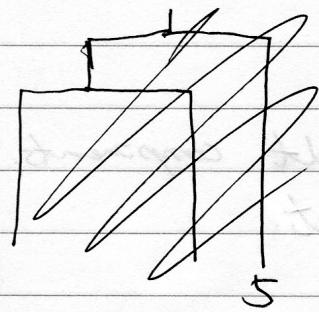
	1	2	3	4
2	5			
3	9	③		
4	10	7	4	
5	16	11	8	6

23	1	23	4
5			
4	10	④	
5	16	8	6

234	5	1	234
		6	

1234 vs 5





Single.

Similarly

# STA 305 Final Web Notes Review

April 12th, 2016

Week 7

- Qualitative Predictors in Regression Models.

- Dummy coding  $(-1, 1, 0)$ .

- $k$  independent categories, then  $k-1$  dummy variables are needed.

$$\text{eg. } X = \begin{cases} 1 & \text{if treatment A} \\ 0 & \text{o.w.} \end{cases}$$

- dummy coding compares each level to the reference level.

- The intercept is the mean of the reference group.

- Deviation coding

- coding system compares the mean of the dependent variable for a given level to the overall mean of the dependent variable.

$$\text{eg. } X_1 = \begin{cases} 1 & \text{green} \\ -1 & \text{yellow} \\ 0 & \text{o.w.} \end{cases} \quad X_2 = \begin{cases} 1 & \text{pink} \\ -1 & \text{yellow} \\ 0 & \text{o.w.} \end{cases} \quad X_3 = \begin{cases} 1 & \text{purple} \\ -1 & \text{yellow} \\ 0 & \text{o.w.} \end{cases}$$

- Estimating treatment effects using least squares.

$y_{ij}$ :  $i$ th treatment,  $j$ th obs.

$$y_{ij} = \mu + T_1 X_{i1} + T_2 X_{i2} + T_3 X_{i3} + \epsilon_{ij}, \quad \epsilon_{ij} \sim N(0, \sigma^2)$$

$$E(y_{Aj}) = \mu_A = \mu$$

$$E(y_{Bj}) = \mu_B = \mu_A + T_1 \Rightarrow T_1 = \mu_B - \mu_A$$

$$E(y_{Cj}) \Rightarrow T_2 = \mu_C - \mu_A$$

$$E(y_{Dj}) \Rightarrow T_3 = \mu_D - \mu_A$$

least square est.

$$\hat{\mu} = \bar{Y}_1$$

$$\hat{T}_1 = \bar{Y}_2 - \bar{Y}_1$$

$$\hat{T}_2 =$$

$$\hat{T}_3 =$$

## Estimating parameters using least squares

$$\left. \begin{aligned} (\text{deviation coding}): E(y_{Aj}) &= \mu_A = T_0 + T_1 \\ E(y_{Bj}) &= \mu_B = T_0 + T_2 \\ E(y_{Cj}) &= \mu_C = T_0 + T_3 \\ E(y_{Dj}) &= \mu_D = T_0 - T_1 - T_2 - T_3 \end{aligned} \right\} \Rightarrow \begin{aligned} T_0 &=? \\ T_1 &=? \\ T_2 &=? \\ T_3 &=? \end{aligned}$$

## Multiple comparisons

Suppose  $k-3$  separate (independent) hypothesis tests conducted

$H_0$  true,  $P(\text{reject } H_0) = \alpha = 1 - P(\text{do not reject } H_0) = 1 - \alpha$

$H_0$  is true,  $P(\text{reject at least one } H_0) = 1 - (1 - \alpha)^k$   
family-wise error rate.

But pairwise error rate is  $P(\text{reject } H_{0,k}) = \alpha$  for any  $k$

- When groups sign. diff from ANOVA:
  - testing all pairs  $\rightarrow$  type I error rate
  - beyond (0.05) that a significant difference is detected when the truth is that no difference exists.

### Bonferroni Method

two sample t-test.  $t_{ij} = \frac{\bar{Y}_{j\cdot} - \bar{Y}_{i\cdot}}{\hat{\sigma} \sqrt{\frac{1}{n_j} + \frac{1}{n_i}}}$   $\hat{\sigma} = \sqrt{MSE}$  from ANOVA table

- trt i & j sig. diff if  $|t_{ij}| > t_{N-k, \alpha/2}$
  - A  $100(1-\alpha)\%$  simultaneous CI for  $c$  pairs  $\mu_i - \mu_j$  is
- $$\bar{Y}_{j\cdot} - \bar{Y}_{i\cdot} \pm t_{N-k, \alpha/2c} \hat{\sigma} \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}$$

### Tukey Method

- The only diff (between T & B) is choice of critical value.

sig diff if  $|t_{ij}| > \sqrt{2} q_{\alpha/2, N-k, \alpha}$

upper  $\alpha$  percentile of the Studentized range dist'n with parameters  $k$  &  $N-k$  degrees of freedom

$$CI: \bar{Y}_j - \bar{Y}_i \pm \frac{1}{\sqrt{2}} q_{\alpha/2, N-k, \alpha} \hat{\sigma} \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}$$

- B is more conservative than T, so CI of T is shorter.

### Sample size for ANOVA

test rejects at  $\alpha$  level if  $\frac{MS_{\text{Treat}}}{MS_E} \geq F_{k-1, N-k, \alpha}$   
 power  $1 - \beta = P(MS_{\text{Treat}}/MS_E \geq F_{k-1, N-k, \alpha})$

- Interpret parameters in the additive model for ANOVA  $\hat{Y}_{ti} = \mu + \cancel{\tau_t} + \epsilon_{ti}$

$$\hat{\mu} = \bar{Y}_1, \hat{\mu} + \hat{\tau}_t = \bar{Y}_t, t=1,2,3,4$$

$$\bar{Y}_t = \sum_{k=1}^{n_t} Y_{tk} / n_t$$

$$\therefore \text{least square est. } \hat{\tau}_t = \bar{Y}_t - \bar{Y}_1, t=1,2,3,4 \quad \epsilon_{ti} \sim N(0, \sigma^2)$$

deviation produced by  $\tau_t$  at  $t$ .

linear model for factorial design:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \beta_4 X_{i4} + \dots + \epsilon_i$$

## Week 8

### factorial design

- fixed # of levels of each of a number of factors
- $2^3$ : 2 levels, 3 factors - 8 experimental conditions
- can involve ~~diff~~ factors with different # of levels.
- ANOVA compares the individual experimental conditions with each other
- factorial --- generally compares combinations of experimental conditions.
- ~~main effects~~
- cube plots
- factorial effects
  - main effects (average effect of a var)
  - interaction effects

$$= \bar{Y}_+ - \bar{Y}_-$$

- replication in factorial design
  - not always feasible actually
  - variance of each measurement is  $\sigma^2$ .
  - estimated variance at each set of condition is

$$S_i^2 = \frac{(Y_{i1} - Y_{i2})^2}{2}$$

$$\text{Pooled est of } \sigma^2 \text{ is } S^2 = \frac{\sum_{i=1}^8 S_i^2}{8} = \frac{\sum \text{diff}^2}{\#}$$

- Estimate of the error variance & standard error of effects from replicated runs

$$\text{Var}(\text{effect}) = \left(\frac{1}{f} + \frac{1}{g}\right) S^2$$

$$SE = \sqrt{\text{Var}}$$

### Interpretation:

interaction plot (parallel = no interaction)

### Linear model for factorial design

\*: The estimated least square coefficients are one-half the factorial estimates, the ~~sample~~ intercept  $\beta_0$  is the sample mean.

### Advantages of factorial design over one-at-a-time designs

- if effect is the same, factorial more efficient, requires fewer obs to achieve the same precision.
- if different, factorial can detect & estimate interactions.

## Week 9

### Randomized Block design

- Comparison of 2 trt:

- unblocked arrangements: unpaired comparison of 2 trt groups.

- blocked arrangements: paired comparison of 2 trts

- Comparison of more trts

- unblocked : randomized one-way design

- blocked : randomized block design

### ANOVA for RB design

~~say k=1, 2, 3~~

ANOVA identity:

- a is the number of trts & b # of blocks

- $y_{i\cdot}$ : total of all obs. taken under treatment i.

- $y_{\cdot j}$ : total of all obs. taken under block j.

- $y_{\cdot \cdot}$ : grand total of all obs.

$N = ab$

$$\bar{y}_{i\cdot} = y_{i\cdot}/b$$

$$\bar{y}_{\cdot j} = y_{\cdot j}/a$$

$$\bar{y}_{\cdot \cdot} = y_{\cdot \cdot}/N$$

$$SS_T = SS_{Treat} + SS_{Block} + SS_E$$

$$\sum_{i=1}^a \sum_{j=1}^b (y_{ij} - \bar{y}_{\cdot \cdot})^2 = b \sum_{i=1}^a (\bar{y}_{i\cdot} - \bar{y}_{\cdot \cdot})^2 + a \sum_{j=1}^b (\bar{y}_{\cdot j} - \bar{y}_{\cdot \cdot})^2 + \sum_{i=1}^a \sum_{j=1}^b (y_{ij} - \bar{y}_{i\cdot} - \bar{y}_{\cdot j} + \bar{y}_{\cdot \cdot})^2$$

- $N$  obs.  $SS_T$  has  $N-1$  dof

- a trts & b blocks so  $SS_{Treat}$  has  $a-1$  &  $SS_{Block}$  has  $b-1$  dof

- $SS_E$  has  $(a-1)(b-1)$  dof

- linear model for RBD:  $y_{ij} = \mu + \tau_i + \beta_j + \epsilon_{ij}$

$$E(\epsilon_{ij}) = 0$$

- Latin Square design: { the design allows for blocking with 2 vars  
the 2 blocking vars must have the same # of levels

- LSD assumes additive, so no interaction effects.
- Graeco-Latin sq. design  
(3 blocking vars) ~~123~~  $\alpha, \beta, \gamma, A, B, C, \dots$
- hyper (>> 3 blocking vars)
- a balanced incomplete block design has the property that every pair of treatments occurs together in a block the same # of times.

## Week 10

- Normal Q-Q plot
  - A marked (systematic) deviation of the plot from the straight line would indicate that:
    - no normality
    - variance is not constant

## Normal plot

- When some of the effects are nonzero the corresponding estimated effects will tend to be larger & fall ~~out~~ off the straight line.
- positive ↑ +
- negative ↓ -

## Half-normal plots

- Lenth's method: testing sig. for experiments without var est.  
 $|ME| < |SME|$ ,  $SME \approx 1.645$  significant
- Blocking factorial designs.

- A general approach for arranging  $2^k$  design in  $2^{k-8}$  blocks of size  $2^{k-8}$  is:  
factorial effects  $v_i$ , block vars  $B_q$
- $B_1 = v_1, \dots, B_8 = v_8$
- e.g.  $B_1 = 135, B_2 = 235, B_3 = 1234, B_4 = 12, B_5 = 145$
- $B_1, B_3 = 245$   
 $B_1, B_2, B_3 = 34$

1. Lower order effects are more likely to be important than higher-order effects.
2. Effects of the same order are equally likely to be important.

- ? Less 2-way - better.

- Fractional factorial designs.

A quarter fraction of the full  $2^5$  design  
or

A  $2^{5-2}$  design.

- Main effect E is aliased with BCD interaction.

$$E = BCD \text{ or } I = BCDE$$

- factorial effects are missing for effects that are aliased.

### Week 11

split plot design is like a RBD.

2 sources of variances,  $\sigma_w^2, \sigma_s^2$

wholeplot      subplot.  
(large)      (small)

whole plot factor A with I levels &  
subplot factor B with J levels  
n times replication.

ANOVA:

Source	DF	SS
Whole-plot replication A	n-1	SS <sub>Rep</sub>
replications x A (whole plot error)	(n-1)(I-1)	SS <sub>w</sub>
Subplot B	J-1	SS <sub>B</sub>
AxB	(J-1)(I-1)	SS <sub>AxB</sub>
Subplot error	I(J-1)(n-1)	SS <sub>s</sub>

linear model of split plot:

$$Y_{ijk} = \gamma + I_k + \alpha_i + (I\alpha)_{ki} + \beta_j + (\alpha\beta)_{ij} + (I\beta)_{kj} + (I\alpha\beta)_{ijk} + \epsilon'_{ijk}, \quad \epsilon'_{ijk} \sim N(0, \sigma^2)$$

$i=1, \dots, I$ ,

$I_k$ : effect of  $k$ th replication

$j=1, \dots, J$

$\alpha_i$ :  $i$ th main effect of A

$k=1, \dots, n$

$(I\alpha)_{ki}$ :  $(k, i)$ th interaction effect

b/w replication & A

↓ this is whole plot error

$\beta_j$ :  $j$ th main effect of B

$(\alpha\beta)_{ij}$ :  $(i, j)$ th interaction between A & B

$\sum = \epsilon_{kij}$

$(I\beta)_{kj}$   $(k, j)$ th

$(I\alpha\beta)_{kij}$

$\epsilon'_{ijk}$ : error term

↓  
subplot error

$\epsilon_{kij} < (I\alpha)_{kj}$  b/c subplot more homogeneous than whole plots.

• Some from Reviews

• Lent's plt. use ME, no ~~ME~~ SME unless asked to adjust for multiple comparison.

• If the design is replicated then it's possible to calc standard errors of factorial estimates.

So we can use sig. testing & CI to evaluate the sig. of effects

So a graphical aid to assess significance is not required.

in conjunction

• If the interaction is sig. then the main effect should be interpreted <sup>in conjunction</sup> with the interaction