

## Index of Lecture 13

Page	Title
1	Practical information
2	Intro sample size issues
3	Statistical methods for sample size
4	Sample size based on estimation precision
5	Errors of type I–II and power
6	Sample size based on power or effect size
7	Non-central distributions
8	Two-way ANOVA with no interaction
9	Calculations for 2-way ANOVA example
10	Unequal sample sizes
11	Sample size misconceptions, and equivalence testing

## PRACTICAL INFORMATION

### Today's lecture:

- “advanced” methods for repeated measures (from L12),
- power and sample size for *experimental studies* — partly well-known material from VHM 801, plus:
  - \* two-way ANOVA with/without interaction
  - \* other advanced methods over VHM 801 in the lab problems,
  - \* software demonstrations (Minitab and Stata) using interactive menus (so no do-file!).

### Reading:

- GO Chapter 7<sup>1</sup> and Section 10.3 (brief),
- supplementary “Notes on sample size calculations” (not part of course curriculum because material mostly covered in textbook).

### Schedule:

- 5th home assignment returned today (with solution),
- last lecture (April 10): project presentations, course wrap-up and *course evaluation*,
- exam date: does morning of April 22 work for everyone?

---

<sup>1</sup> Discussion about non-central  $F$ -distributions and power curves may be skipped because our focus is on using software for power calculations.

## INTRO SAMPLE SIZE ISSUES

Factors affecting the size of an experimental study:

- # treatments to compare (incl. control),
- # factors and for each factor, the number of levels,
- experimental design, e.g. block sizes in a block design,
- cost per experimental unit,
- management of experiment (e.g., size limitations).

Statistical considerations:

- size should be sufficient to detect (statistical significance of) treatment differences of interest,
- avoid “waste” of experimental units,
- reduce sensitivity to errors (by taking replications).

Textbook example (GO 7.1-5, p. 150 ff.):

- patients with 3 different types of neurological diseases,
- “VOR” measurements as indicators of disease status,
- preliminary data (log scale): group means 2.82, 3.89 and 3.04, within-group variance 0.075.

Additional notes example:

- within-subject differences<sup>2</sup> in blood pressure with an assumed standard deviation of 10 *mm* Hg.

---

<sup>2</sup> E.g., differences computed from measurements before and after an intervention.

## STATISTICAL METHODS TO CHOOSE SAMPLE SIZE

Fact: all procedures require pre-decided statistical model and detailed prior knowledge (estimates or guesses) about the outcomes:

- size of effects or precision of interest,
- standard deviation of observations (for normal data<sup>3</sup>).

Overview of approaches for determining sample size:

- from desired precision (standard error, size of 95% CI) on selected estimate (treatment mean, contrast),
- from effect of interest, using “Cox’s rule” (practical rule),
- from desired power of test for effect of interest, preferably using statistical software:
  - \* Minitab/Stata (or websites) for basic designs,
  - \* SAS version 9 for a range of complex designs,
  - \* specialised software for special/advanced designs<sup>4</sup>:
    - (active researcher): <http://www.stat.uiowa.edu/~rlenth/Power>
    - (stats webpages): <http://statpages.org/#Power>
    - (G\*Power = free power calculation software): <http://www.gpower.hhu.de/>

or using *simulation* (some refs in the notes).

<sup>3</sup> For data with the variance estimated independently of the mean

<sup>4</sup> Check also Stata’s add-on packages (search sample size), and the UCLA intro webpage [http://www.ats.ucla.edu/stat/seminars/Intro\\_power/default.htm](http://www.ats.ucla.edu/stat/seminars/Intro_power/default.htm) and the accompanying examples on <http://www.ats.ucla.edu/stat/dae/>.

## SAMPLE SIZE BASED ON ESTIMATION PRECISION

Normal distribution models:

- assume error std.dev.  $\sigma$  and estimate/guess of value,
- assume parameter of interest  $Par$  and estimate  $Est$ , with standard error  $SE(Est) = \sigma A$ , where  $A = A(n)$  is a known constant (depending on number of obs.  $n$ ),
- approximate<sup>5</sup> 95% CI:  $Est \pm 2\sigma A(n)$ .

Compute  $n$  to achieve desired margin of error (or CI length) by solving with respect to  $n$  in the equation:

$$\text{desired value} \geq 2\sigma A(n).$$

Blood pressure example: desired margin of error: 3 mm Hg,

- model:  $n$  i.i.d. observations (differences!) from  $N(\mu, \sigma^2)$ ,
- $Par$  = population mean,  $Est$  = sample mean,  $A = 1/\sqrt{n}$ ,
- solve:  $3 \geq 2 \times 10/\sqrt{n} \Rightarrow n \geq (2 \times 10/3)^2 = 44.4 \approx 45$ .

VOR example: desired CI of length 0.5 for mean differences between (any) two groups,

- model: 3 independent samples of size  $n$  from normal distributions  $N(\mu_i, \sigma^2)$ ; use  $\sigma = s_p = \sqrt{0.075} = 0.2739$ ,
- $Par = \mu_1 - \mu_2$ ,  $Est = \bar{y}_1 - \bar{y}_2$ ,  $A = \sqrt{2/n}$ ,
- solve:  $0.25 \geq 2 \times 0.2739 \sqrt{2/n} \Rightarrow n \geq 2(2 \times 0.2739/0.25)^2 = 9.6$ , or  $n = 10$ ; rerun with 2 replaced by  $t(.975, 27) = 2.052$ ;  $n \geq 10.1$ ; choose  $n = 11$  patients per group.

---

<sup>5</sup> Valid for  $\sigma$  unknown and df large (say  $\geq 40$ ), so that  $t_{.975}(\text{df}) \approx z_{.975} \approx 2$ .

## ERRORS OF TYPE I–II AND POWER

Errors of type I and II:

- type I error: to reject  $H_0$ , when  $H_0$  in reality true,
- type II error: to not reject  $H_0$ , when  $H_0$  in reality false,
- definition of statistical tests involves only type I errors, which are controlled by the significance level ( $\alpha$ ),
- power of statistical tests involves type II errors (below),
- overview:

Conclusion from sample	Truth about population	
	$H_0$ true	$H_a$ true ( $H_0$ false)
reject $H_0$	type I error	no error
not reject $H_0$	no error	type II error

Power of a statistical test:

- involves a specific alternative, e.g. in the blood pressure example a difference of 3 *mm* Hg,
- definition: power = probability that the statistical test will *reject*  $H_0$ , when the specific alternative  $H_a$  is true,  
= 1 – type II error,
- important for planning of experiments: what chance of a significant result?
- difficult to calculate in all models (use software!),
- typical values for planning of a study: 0.8, 0.9, 0.95.

## SAMPLE SIZE BASED ON POWER OR EFFECT SIZE

Cox's rule for “informal” sample size determination:<sup>6</sup>

- assume effect size  $|m_0 - m_a|$ , where  $m_0$  and  $m_a$  are (true, population) values of *Par* under  $H_0$  and  $H_a$ , respect.,
- rule: choose  $n$  by solving:  $|m_0 - m_a| = 3\sigma A(n)$ ,
- note: does not involve power or significance level.

Requirements for sample size calc. based on power:<sup>7</sup>

- statistical model / design and corresponding software,
- size of effect desired to be detected,
- standard deviation of model (normal data),
- desired value of power (0.8, or 80%, commonly used),
- signif. level and direction (one/two-sided) of test used.

Blood pressure example: true difference of 3 *mm* Hg,

- Cox's rule:  $3 = 3 \cdot 10\sqrt{1/n} \Rightarrow n = 100$ ,
- Minitab with power 0.8 and sign. level 0.05:  $n = 90$ .

VOR example: power for  $n = 4$  and sign. level 0.01 with means as observed in preliminary data:

- Textbook/Russell Lenth software: 0.930 (all means),
- Minitab (vers. 16) / Stata: 0.896 (maximal difference).

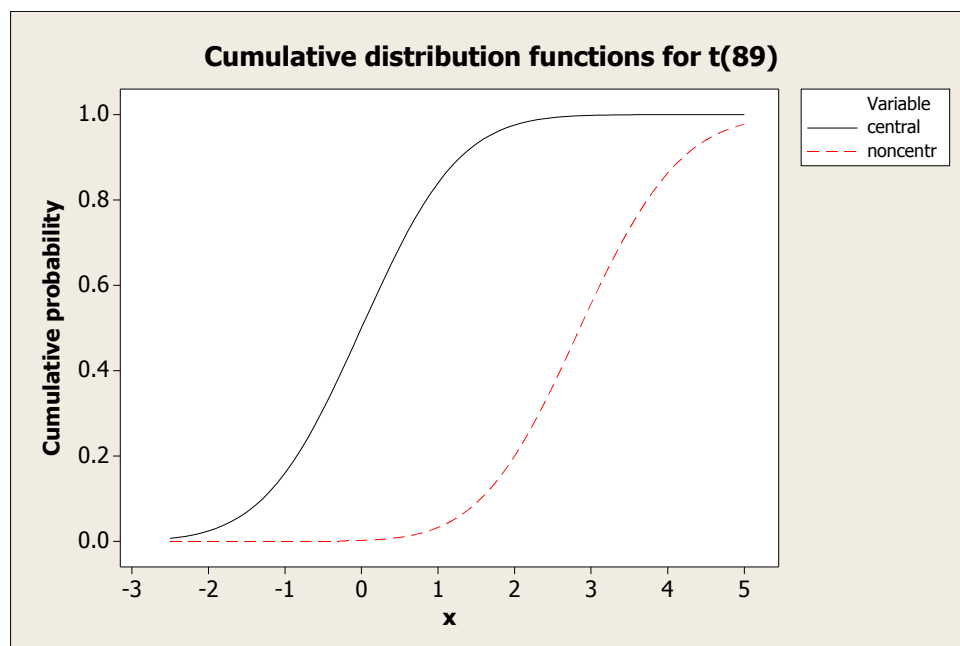
<sup>6</sup> Discussed in Christensen (1996): *Analysis of Variance, Design, and Regression*.

<sup>7</sup> Post-hoc power calculation (after study has been carried out) is controversial and *not* recommended.

## NON-CENTRAL DISTRIBUTIONS

Several of the common ref. distrib. for tests ( $t$ -,  $\chi^2$ -,  $F$ -distrib.) have a non-centrality parameter ( $\zeta$ ):

- $\zeta = 0 \sim$  usual ref. distrib. under null hypothesis  $H_0$ ,
- $\zeta \neq 0 \sim$  distrib. of test statistic under specific alternative hypothesis  $H_a$ , where  $\zeta \sim$  deviation from  $H_0$ ,
- power is computed as tail areas in distrib. with  $\zeta \neq 0$ ,
- example: non-central  $F$ -distribution in ANOVA table – see GO Figure 7.1, where  $\zeta = \sum_i n_i \alpha_i^2 / \sigma^2$ ,
- blood pressure example: non-central  $t$ -distribution:
  - \*  $t = \hat{\mu} / \text{SE}(\hat{\mu}) \sim$  non-central  $t$ , where  $\hat{\mu} \sim N(\delta, \sigma_{\hat{\mu}}^2)$  and  $\zeta = \delta / \sigma_{\hat{\mu}}$ ,
  - \* with  $n = 90$  we get  $\zeta = 3 / (10 / \sqrt{90}) = 2.846$  and  $t^* = t_{.975, 89} = 1.987$ :



## SAMPLE SIZE FOR TWO-WAY ANOVA

Example 13.5-6 in IPS (5th/6th/7th ed.): calcium supplementation as a treatment for Osteoporosis in elderly people,

- factor: calcium (placebo, 800 *mg*/day),
- factor: vitamin D (placebo, 300 IU/day)<sup>8</sup>,
- outcome: change in bone mineral density (BMD),
- assume  $n$  subjects per calcium  $\times$  vitamin D group,<sup>9</sup>
- assume calcium effect size of interest:  $\delta = 5$  BMD units,
- assume within-group standard deviation, *for the change in BMD*:  $\sigma = 10$  units.

Statistical model (anticipated):

$$Y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk}$$

where  $\alpha_i$ ,  $\beta_j$ ,  $(\alpha\beta)_{ij}$  are the main effects and interaction, respectively, and  $\varepsilon_{ijk} \sim N(0, \sigma^2)$ .

Two possible situations:

- no interaction: sample size based on main effects,
- interaction, and different approaches for sample size:
  - \* based on contrast of interest for one factor within one level of other factor ( $\approx$  two-sample situation),
  - \* based on  $F$ -test for interaction,
  - \* based on contrast within interaction (for  $2 \times 2$  factorial: equivalent to full interaction).

---

<sup>8</sup> Vitamin D is needed for the body to efficiently utilize calcium.

<sup>9</sup> In the notes, the number of subjects per group is denoted by  $c$ .

## CALCULATIONS FOR 2-WAY ANOVA EXAMPLE

Sample size for calcium effect ( $\delta$ ) in no interaction model:

- two-sample calc. (power 0.80, Minitab) gives  $n = 64$   
 $\Rightarrow$  actual  $n = 64/2 = 32$ , because the two vitamin D groups both contribute to the sample size,<sup>10</sup>
- using Cox's rule:  $SE = \sigma \sqrt{1/(2n) + 1/(2n)} = \delta/3$   
 $\Rightarrow n = (3\sigma/\delta)^2 = 36$ .

Interaction model – two approaches:

- sample size for calcium effect in high (or low) vitamin D group: above two-sample calc. applies<sup>9</sup>:  $n = 64$ ,
- sample size for interaction of size  $\delta$ :
  - \* interaction contrast<sup>11</sup> and SE:

$$\begin{aligned}\hat{\gamma} &= \bar{Y}_{11\cdot} + \bar{Y}_{22\cdot} - \bar{Y}_{12\cdot} - \bar{Y}_{21\cdot}, \\ SE(\hat{\gamma}) &= \sigma \sqrt{1/n + 1/n + 1/n + 1/n} = \\ &= \sigma \sqrt{4/n} = (\sqrt{2}\sigma) \sqrt{2/n},\end{aligned}$$

- \* last formula  $\sim$  2-sample comparison with standard deviation  $\sqrt{2} \cdot \sigma$ :  
 $\Rightarrow n = 127$  (power 0.80, Minitab).

<sup>10</sup> Actual power will be slightly less than 0.80 because the df is smaller in the two-way ANOVA than the two-sample situation.

<sup>11</sup> The contrast estimates the difference in calcium effect between the two vitamin D levels, or the difference in vitamin D effect between the two calcium levels.

## UNEQUAL SAMPLE SIZES

Generally speaking, it is most efficient (gives best precision) to have equal sample sizes in different groups of a factor.

Two special cases where unequal distributions across groups are useful:

- One control group and several treatment groups all to be compared to control only(!): the optimal sample size for control is larger than for treatment groups,
  - \*  $g$  treatments (incl. control),  $n_c$ =size of control,  $n_t$ =size for other treatments  $\Rightarrow$   
solve equation:  $(n_c/n_t)^2 = (g - 1)$ ,
  - \* e.g., for  $g = 5$  we get  $n_c = 2n_t$  (control group of double size!),
- factor with quantitative (and equidistant) levels where certain polynomial contrasts are of particular interest:
  - \* e.g., linear contrasts always give highest weights to the most extreme categories (Table D.6 in GO)  $\Rightarrow$  higher precision for linear contrast may be achieved by overrepresenting the most extreme categories.<sup>12</sup>

---

<sup>12</sup> The gain in precision is counterbalanced by reduced precision for other contrasts; for example, if only the two extreme categories are included, the linear contrast is estimated with best precision, but no other polynomial contrasts can be estimated.

## SAMPLE SIZE MISCONCEPTIONS, AND EQUIVALENCE TESTING

Common misconceptions<sup>13</sup> in sample size calculations:

- use of standard effect sizes (general definitions of “small”, “medium” and “large” effects, relative to std. dev.): effects of interest should be determined exclusively from the context of your study,
- retrospective power calculation: after a study has been carried out and using its estimated values:
  - \* power/sample size calculations aid in planning of new studies, not in interpreting results of data analysis,
  - \* confidence intervals give the best information about the unknown parameters from a study,
  - \* if  $H_0$  was not rejected, the conclusion may be strengthened by an equivalence test (instead of arguing from the study’s power).

An equivalence test is for making a statement that effects of two “treatments” differ at most by a (biologically) small amount (say  $\delta$ ),

- for  $H_0 : \theta = 0$  ( $\theta$  being a difference in means, or other parameters), not rejecting  $H_0$  is a weak and non-quantitative conclusion,
- a CI for the difference  $\theta$  contains useful information,
- a non-equivalence hypothesis  $H_0^{(ne)} : |\theta| \geq \delta$  can be tested against the  $H_a^{(ne)} : |\theta| < \delta$ , as follows (at a 5% significance level):
  - \* compute a 90% CI (not 95% CI) for  $\theta$ ,
  - \* reject  $H_0^{(ne)}$ , if the interval  $(-\delta, \delta)$  is entirely inside the CI.

---

<sup>13</sup> Largely based on Lenth (2001), *The American Statistician* **55**, 187–193.