

## **Experimental design**

**Randomization:** assigning experimental units to treatments at random

- Eliminates conscious or unconscious bias in assigning units to treatments
- Spreads variability

**Random sample:** each of  $N$  items is equally likely to appear in a sample of size  $n$

- Observational study: permits valid inferences about population parameters
- Experimental study: that (above), plus serves as a method of randomization (assign  $n_i$  units to treatment  $i$ )

## Randomization algorithms

### A. Shuffling method

0. Label the units 1 to  $N$
1. Generate  $N$  uniform(0, 1) random variables and associate each one with the unit labels

$$\begin{array}{ll} 1 & u_1 \\ 2 & u_2 \\ 3 & u_3 \\ & \vdots \\ N & u_N \end{array}$$

2. Sort the observations by the uniform variable
3. Pick the first  $n$  observations in the sorted data

### B. Fast method (one pass through the data)

0. Read the next unit label and make it current
1. Generate  $U \sim \text{uniform}(0, 1)$
2. Test if  $NU > n$ . If yes, go to step 6
  3. Include current label in sample
  4. Set  $N = N - 1$  and  $n = n - 1$
  5. If  $n > 0$  go to step 1; otherwise terminate
6. Skip over current label
  7. Set  $N = N - 1$  and go to step 1

(Vitter, J. S. 1984. Faster methods for random sampling. Communications of the ACM 27:703-718)

## How big a sample?

Question revolves around:

- how precise an estimate is desired by investigator
- how big a departure from  $H_0$  that the investigator desires to have a good chance of detecting

Estimate precision:  $100(1 - \alpha)\%$  CI is in the form

$$\hat{\theta} \pm \hat{E}$$

where  $\hat{E}$  is the **margin of error** ( $\frac{1}{2}$ -width of CI). Strategy: fix size of  $\hat{E}$  and solve for the value of  $n$  that “makes it so”.

**ex.** One-sample estimate of  $\mu$  in a normal( $\mu, \sigma^2$ ) population

$$\bar{x} \pm t_{\alpha/2} \sqrt{\hat{\sigma}^2/n}$$

Note:  $t_{\alpha/2}$  depends on  $n$  ( $n - 1$  df)  
 $\sigma^2$  unknown (and  $\hat{\sigma}^2$  will not result until the sample is drawn)

So: approximate  $t_{\alpha/2}$  with  $z_{\alpha/2}$  (large  $n$ )  
guess  $\sigma^2$  (preliminary study?)

Then:  $E = z_{\alpha/2} \sqrt{\sigma^2/n} \Rightarrow$  pick  $n = \frac{z_{\alpha/2}^2 \sigma^2}{E^2}$

**ex.** Estimating  $p$  in a binomial( $n, p$ ) sample

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

Guess  $p$  (or take  $p = \frac{1}{2}$  as worst case)

$$E = z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \Rightarrow \text{pick } n = \frac{z_{\alpha/2}^2 p(1-p)}{E^2}$$

**Hypothesis tests:** typically in the form

$$H_0: \theta = \theta_0 \quad (\theta \text{ a parameter})$$

$$H_a: \theta \neq \theta_0$$

Test statistic  $S$ : reject  $H_0$  if  $S \geq c$

Under  $H_0$ ,  $S$  has a  $\left\{ \begin{array}{l} \text{t distribution} \\ \text{F distribution} \\ \text{chi-square dist} \end{array} \right.$  etc. depending on

application. But what is the distribution of  $S$  if  $H_a$  is true?

$S$  has a  $\left\{ \begin{array}{l} \text{noncentral t distribution} \\ \text{noncentral F distribution} \\ \text{noncentral chi-square dist} \end{array} \right.$

Noncentral distributions depend on:

sample size

$|\theta - \theta_0|$ , the effect size

Design strategy: fix the effect size one wants to be able to detect, and fix  $1 - \beta$  (the power of the test); solve for the sample size that “makes it so”.

In SAS: inverse distribution functions (for calculating critical values) and noncentral distributions are library functions

TINV( $pr, df, \lambda$ )	PROBT( $x, df, \lambda$ )
FINV( $pr, df_1, df_2, \lambda$ )	PROBF( $x, df_1, df_2, \lambda$ )
CINV( $pr, df, \lambda$ )	PROBCHI( $x, df, \lambda$ )

Here,  $\lambda$  is the “noncentrality parameter” which is related to the effect size (formula for  $\lambda$  varies between applications).

**ex.** One-way AOV:  $Y_{ij} \sim \text{normal}(\mu_i, \sigma^2)$   $i = 1, 2, \dots, t$   
 $j = 1, 2, \dots, n_j$   $n_1 + n_2 + \dots + n_t = n_T$

The noncentrality parameter is

$$\lambda = \frac{\sum_{i=1}^t n_i (\mu_i - \mu_{.})^2}{\sigma^2}$$

If the smallest distance between means that investigator wants to detect is  $D = |\mu_k - \mu_l|$ , and investigator will assign  $n_i = r$  units to every treatment, then the minimum value of  $\lambda$  is

$$\lambda = \frac{rD^2}{2\sigma^2}$$

Strategy: fix  $D$ , guess  $\sigma^2$  (preliminary study, published info, etc), and then use the infamous Pearson-Hartley charts (yuck) or the handy SAS program (woo hoo!) to find the  $r$  value that “makes it so”.