SAMPLE SIZE CALCULATIONS

Determining how large a sample should be drawn is a big part of practical statistics.

More samples amount to more information, which translates into navrower confidence intervals for unknown population parameters.

In these notes, we derive formulas for the sample sizes needed to achieve a desired narrowness of (1-d) 100% C.I.s for the population mean μ , the population proportion p, and for a difference in population means $\mu_1 - \mu_2$ or proportions $p_1 - p_2$.

ONE-SAMPLE SETTING

For a mean M:

Let $X_{1,...,}X_n$ be ind rus with mean μ , variance σ^2 . Then a large-n $(1-d)^4/00\%$ C.I. for μ has the form

$$\bar{X}_n \stackrel{f}{=} E_{\alpha/2} \stackrel{f}{\int n} M.E.$$

The part added and subtracted is called the <u>margin of error</u> (M.E.). The M.E. decreases as the sample size n zrows, making the C.I. narrower.

We approach sample size calculation as follows: Choose an upper bound Mt for the M.E. and find the smellest n such that M.E. $\leq M^*$.

The value of M⁴ comes from the scientific researcher, and it reflects the degree of accuracy with which it is desired to estimate the unknown parameter.

where FX7 is the smallest integer greater than or equal to x.

- To make this calculation in practice, we replace σ^2 with an estimate $\hat{\sigma}^2$ from a previous or "pilot" study.
- <u>Example</u>: Researchers want to build a 95% C.T. for n with a M.E. no greater than $\frac{1}{2}$. It is believed that $\sigma \approx 2$. What sample size n should they take?

Take

$$n = \left[\left(\frac{\overline{1.96}}{V_2} \right)^2 (2)^2 \right] = \left[61.4656 \right] = 62.$$

 M^*

For a proportion p:

Let
$$X_{i,j,...,j} X_n \stackrel{iid}{\sim} Bernoulli(p)$$
.
Then a large-n $(1-d)^{\pm} 100 \%$ C.I. for p has the form
 $\hat{p}_n \pm Z_{oly} \int \frac{p(1-p)}{n}$
M.E.

To achieve M.F. = M, we choose the smallest n such that

$$Z_{a_{2}}\sqrt{\frac{p(1-p)}{n}} \leq M^{*} \quad \langle z \rangle \quad n \geq \left(\frac{Z_{a_{2}}}{M^{*}}\right)^{p(1-p)},$$

i.e., we choose
$$n = \left\lceil \left(\frac{2}{M^2}\right)^2 \neq (1-\frac{1}{2}) \right\rceil$$
.

In practice, we replace & with either

- (i) an estimate \$ from a previous or pilot study.
- (ii) the value $p = \frac{1}{2}$, for which the M.E. is maximized at any fixed sample size n. Using $p = \frac{1}{2}$ in our sample size calculation gives the most conservative (largest) value of n.
- Example: Suppose we wish to build a 91% C.I. for the proportion of voters who will select a certain candidate. We went a M.E. of no more than 2 percentage points. What sample size do we need?

Take
Take

$$n = \left[\left(\frac{2.576}{.02} \right)^2 + \frac{1}{2} (1-\frac{1}{2}) \right] = \left[\frac{1}{2} \frac{1}{$$

TWO-SAMPLE SETTING

For comparing means:
For large
$$n_1, n_2$$
, the interval
 $\overline{X} - \overline{Y} \pm \frac{1}{2} \frac{\sigma_1^2}{\sigma_1} + \frac{\sigma_2^2}{\sigma_2}$
is an approximate $(1 - d)^2 100\%$ C.T. for $p_1 - p_2$.
Let $n = n_1 + n_2$ be the total sample size and let
 $d = n_1/n$
 $1 - \delta = \frac{n_2}{n}$,

3

So that $T \in (0,1)$ is the proportion of observations coming from population 1.

Then the M.E. for a given n and de(0,1) can be written

$$\mathcal{M}(n, \delta) = Z_{a_{2}} \sqrt{\frac{\sigma_{1}^{2}}{\delta n} + \frac{\sigma_{2}^{2}}{(1-\delta)n}}$$

We observe

(i) M(n, t) is strictly decreasing in n. (ii) For each n, $M(n, \cdot)$ is minimized at $t_{opt} = \frac{\sigma_1}{\sigma_1 + \sigma_2}$. (iii) For each n, $M(n, \cdot)$ has the minimum $Z_{al_2}(\frac{\sigma_1 + \sigma_2}{\sqrt{n}})$, i.e.

$$M(n, t_{opt}) = Z_{a/2} \left(\frac{\sigma_1 + \sigma_2}{\sqrt{n}} \right)$$

To achieve a M.E. = M^+ , we choose n_1 and n_2 as follows:

• Find smallest real number nt such that

$$\frac{z}{d_{2}} \frac{(\sigma_{1} + \sigma_{2})}{\sqrt{n^{4}}} \leq M^{4},$$

i.e. get
$$n = \left(\frac{\overline{z}}{M^{4}}\right)^{2} \left(\sigma_{1} + \sigma_{2}\right)^{2}.$$

Then set $n_{1} = \left[\left(\frac{\sigma_{1}}{\sigma_{1} + \sigma_{2}}\right)^{n^{4}}\right], \quad n_{2} = \left[\left(\frac{\sigma_{2}}{\sigma_{1} + \sigma_{2}}\right)^{n^{4}}\right], \quad n = n_{1} + n_{2}.$

In practice, we replace σ_1 and σ_2 with some estimates $\hat{\sigma}_1$ and $\hat{\sigma}_2$ from a previous or pilot study.

4

Note that we draw a larger sample from the population with the larger variance.

<u>To find lopt</u>: Find value of I which minimizes $M(n, \delta)$. Since $M(n, \delta) \ge 0$, the same value of I minimizes $M^2(n, \delta)$. So we get the derivative of $M^2(n, \delta)$ w.r.t. I equal to zero and solve for I:

$$\frac{2}{24} M^{2}(n, \delta) = \frac{2}{24} \frac{2}{54} \left[\frac{\sigma_{1}^{2}}{4n} + \frac{\sigma_{2}^{2}}{(1-\delta)^{n}} \right] = \frac{2}{54} \frac{2}{54} \left[-\frac{\sigma_{1}^{2}}{5^{2}} + \frac{\sigma_{2}^{2}}{(1-\delta)^{2}} \right] = 0$$

$$(2)$$

$$\frac{\sigma_{2}^{2}}{(1-\delta)^{2}} = \frac{\sigma_{1}^{2}}{5^{2}}$$

$$(2)$$

$$\frac{\sigma_{1}^{2}}{(1-\delta)^{2}} = \frac{\sigma_{1}}{5^{2}}$$

$$\frac{\sigma_{1}}{5^{2}} = \frac{\sigma_{1}}{5^{2}}$$

$$(2)$$

$$\frac{\sigma_{1}}{5^{2}} = \frac{\sigma_{1}}{5^{2}}$$

$$\underline{\mathbf{T}_{o_{2}} \mathbf{t}} \underbrace{M\left(n, t_{opt}\right)}_{M\left(n, t_{opt}\right)} : \text{Under } t_{opt}, \text{ we } \mathbf{y}_{opt} \\
M\left(n, t_{opt}\right) = \mathbb{E}_{d_{12}} \sqrt{\frac{\sigma_{1}^{2}}{t_{opt} \cdot n} + \frac{\sigma_{2}^{2}}{\left(1 - \delta_{opt}\right) \cdot n}} \\
= \mathbb{E}_{d_{12}} \sqrt{\frac{\sigma_{1}^{2}}{\left(\frac{\sigma_{1}}{\sigma_{1} + \sigma_{2}}\right)n} + \frac{\sigma_{2}^{2}}{\left(1 - \frac{\sigma_{1}}{\sigma_{1} + \sigma_{2}}\right)n}} \\
= \mathbb{E}_{d_{12}} \sqrt{\frac{\sigma_{1}\left(\sigma_{1} + \sigma_{2}\right) + \sigma_{2}\left(\sigma_{1} + \sigma_{2}\right)}{n}} \\
= \mathbb{E}_{d_{12}} \sqrt{\frac{\sigma_{1}^{2} + 2\sigma_{1}\sigma_{2} + \sigma_{2}^{2}}{\left(\frac{\sigma_{1}^{2} + \sigma_{2}^{2}\right)}}} \\
= \mathbb{E}_{d_{12}} \left(\frac{\sigma_{1} + \sigma_{2}}{\sqrt{n}}\right).$$

5

Example: Suppose we have
$$\hat{\sigma}_1 = 2$$
 and $\hat{\sigma}_2 = 3$ from a previous study.
We wish to build a 97% C.T. for $\mu_1 - \mu_2$ which
has a M.E. ≤ 0.5 . Find n_1 and n_2 .
Get $n_1^2 = \left(\frac{2.596}{0.5}\right)^2 (2+3)^2 = 663.4897$,

and then set

$$n_1 = \left\lceil \left(\frac{2}{2+3}\right) 663.4897 \right\rceil = 266$$

 $n_2 = \left\lceil \left(\frac{3}{2+3}\right) 663.4897 \right\rceil = 399$, $n = 266 + 399 = 665$.

For comparing proportions:

For large
$$n_{1,n_{2,s}}$$
 the interval
 $\hat{p}_{1} - \hat{p}_{2} \pm \Xi_{d/2} \sqrt{\frac{p_{1}(1-p_{1})}{n_{1}} \pm \frac{p_{2}(1-p_{2})}{n_{2}}}$
is an approximate $(1-d)^{4}/00\%$ C.T. for $p_{1}-p_{2}$.
By the same arguments as those in the previous section,
we may justify choosing n_{1} and n_{2} as follows:
• Find smallest real number n^{4} such that
 $\Xi_{d/2} \left(\sqrt{p_{1}(1-p_{1})} + \sqrt{p_{2}(1-p_{2})} \right) \leq M^{4}$
i.e. get $\pm (\Xi_{d/2})^{2} (1-p_{2}) + \sqrt{p_{2}(1-p_{2})}$

i.e.
$$qe^{\frac{1}{2}}$$

 $n^{\frac{1}{2}} = \left(\frac{\frac{1}{2}d_{2}}{M^{\frac{1}{2}}}\right)^{2} \left(\left[\frac{1}{p_{1}}(1-p_{1}) + \frac{1}{p_{2}}(1-p_{2})\right]\right)$.

• Then set
$$n_1 = \left[\left(\frac{\sqrt{p_1(1-p_1)}}{\sqrt{p_1(1-p_1)} + \sqrt{p_2(1-p_2)}} \right) n^4 \right]$$

 $n_2 = \left[\left(\frac{\sqrt{p_2(1-p_2)}}{\sqrt{p_1(1-p_1)} + \sqrt{p_2(1-p_2)}} \right) n^4 \right]$ and $n = n_1 + n_2$.

In practice, we replace β_1 and β_2 with some estimates $\hat{\beta}_1$ and $\hat{\beta}_2$ from a previous or pilot study.