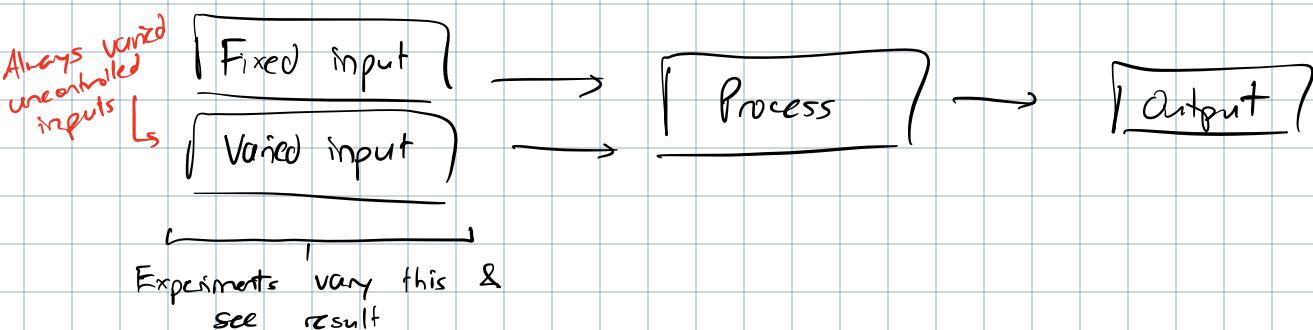
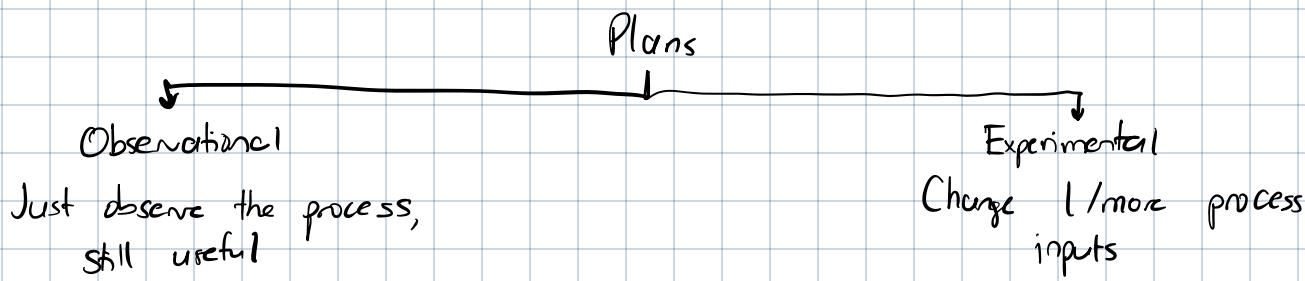


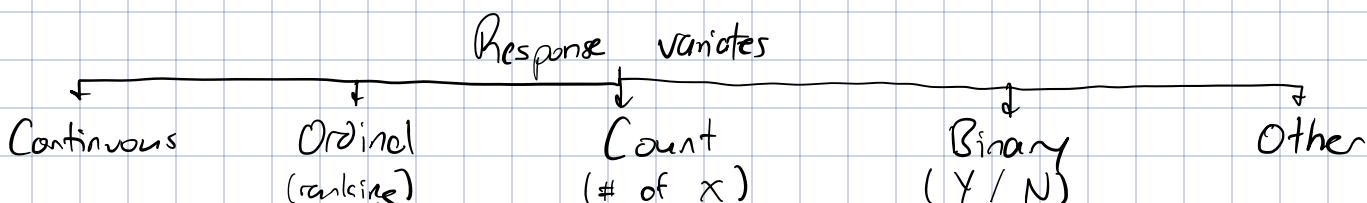
FUNDAMENTALS OF EXPERIMENTAL PLANS



Factor: single variable that is changed / set on each unit in sample, focus of study

↳ Factor levels: set of values assigned to factor (quantitative / qualitative), categories

Treatment: combo of factor levels that can be applied to a unit



Fundamental Principles of Experimental Plans

Blockings: groups (blocks) of units where 1/more explanatory variables are fixed while diff. treatments applied to units w/in group

↳ Ex:// Sample consists of 300 stores. Stores put in blocks of 4 where all stores of 1 block had similar 3-week sales volume. 4 diff. treatments applied
expl. variable held fixed.

↳ Pros:

① Prevents confounding as explanatory variables fixed

② Improves precision of cond. b/c w/in-block variation ↓

Confounding: rel. b/w 2 inputs & vary together ⇒ cannot det. if outcome from 1 or other

Replication: apply treatment to > 1 unit in sample

→ Ex:// Treatment A applied to 7S stores (once in 7S blocks)

↳ Pros:

- ① Reduces sample error b/c larger sample possible
- ② Estimate conclusion precision $\Rightarrow \hat{\sigma}^2 = \frac{\sum (x_i - \bar{x})^2}{n-1} \quad n > 1$!

Random assignment: assign treatments to units randomly

→ Ex:// Each store in block had an equal chance of getting Treatment A

↳ Pros:

- ① Reduces confounding risk
- ② Generates analysis method

INFEERENCE OVERVIEW

Gaussian Model

Model: $y_i = \mu + \epsilon_i, \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$
↳ 'true' avg. of pop.

Assumption checking:

- ① Normality \Rightarrow histogram, Q-Q plot
- ② Constant variance \Rightarrow residual analysis
- ③ ϵ_i are indep. \Rightarrow can't check rely on randomization

t-discrepancy measure

To fit model, estimate μ & σ via max. likelihood (MLE) or least squares (LSE)

LSE of μ & σ = MLE of μ & σ for Gaussian model

↳ To estimate μ :

$$\min_{\mu} \sum_{i=1}^n (y_i - \mu)^2 \Rightarrow \hat{\mu} = \bar{y}$$

↳ To estimate σ :

$$\hat{\sigma} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}$$

Estimator distn. of μ & σ :

$$\tilde{\mu} \sim N(\mu, \frac{\sigma^2}{n})$$

↳ Proof:

(1) Definition of estimator:

$$\tilde{\mu} = \frac{1}{n} \sum_{i=1}^n y_i$$

(2) Normality:

$y_i \sim N(\mu, \sigma^2)$. Since $\tilde{\mu}$ is a linear combo. of y_i , $\tilde{\mu}$ is also normally distributed. indep.

(3) Mean:

$$E\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = \frac{1}{n} \sum_{i=1}^n E(y_i) = \frac{1}{n} \sum_{i=1}^n \mu = \mu$$

(4) Variance:

$$\text{Var}\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(y_i) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}$$

$$\frac{\hat{\sigma}^2(n-1)}{\sigma^2} \sim \chi^2_{n-1}$$

↳ Proof: We know that

$$\sum_{i=1}^n (y_i - \mu)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 + n(\bar{y} - \mu)^2$$

Now:

$$\sum_{i=1}^n \left(\frac{y_i - \mu}{\sigma} \right)^2 = \underbrace{\sum_{i=1}^n \left(\frac{y_i - \bar{y}}{\sigma} \right)^2}_{U} + \underbrace{\left(\frac{\bar{y} - \mu}{\sigma/\sqrt{n}} \right)^2}_{V} + \underbrace{\left(\frac{\bar{y} - \mu}{\sigma/\sqrt{n}} \right)^2}_{W}$$

Since $\frac{y_i - \mu}{\sigma} \sim N(0, 1)$, $U \sim \chi^2_n$.

From above, $\bar{y} \sim N(\mu, \sigma^2/n) \Rightarrow \frac{\bar{y} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \Rightarrow W \sim \chi^2_1$

Furthermore:

$$\tilde{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

Thus:

$$V = \frac{(n-1) \tilde{\sigma}^2}{\sigma^2}$$

Assuming \bar{Y} is indep. of σ^2 , V is indep. of ω . Thus

$$M_V(t) = M_V(t) M_W(t)$$

$$\frac{1}{(1-2t)^{n/2}} = M_V(t) \cdot \frac{1}{(1-2t)^{k/2}}$$

$$M_V(t) = \frac{1}{(1-2t)^{n-1/2}}$$

Thus:

$$V = \frac{(n-1) \tilde{\sigma}^2}{\sigma^2} \sim \chi^2_{n-1}$$

Note that:

$$Z \sim N(0, 1), V \sim \chi^2_k, Z \text{ is indep. of } V \Rightarrow T = \frac{Z}{\sqrt{V/k}} \sim t_k$$

Claim:

$$\frac{\tilde{\mu} - \mu}{\sqrt{\tilde{\sigma}^2/n}} \sim t_{n-1}$$

Let's Proof:

Let

$$Z = \frac{\tilde{\mu} - \mu}{\sqrt{\tilde{\sigma}^2/n}} \sim N(0, 1), V = \frac{(n-1) \tilde{\sigma}^2}{\sigma^2} \sim \chi^2_{n-1}$$

We know that $\tilde{\mu}$ is indep. of $\tilde{\sigma}^2$, so Z & V are indep.
b/c func. of indep. r.v. are also indep.

Thus:

$$\begin{aligned} \frac{Z}{\sqrt{V/k}} &= \frac{\frac{\tilde{\mu} - \mu}{\sqrt{\tilde{\sigma}^2/n}}}{\sqrt{\frac{(n-1)\tilde{\sigma}^2}{\sigma^2} \cdot \frac{1}{n-1}}} \\ &= \frac{\tilde{\mu} - \mu}{\sqrt{\tilde{\sigma}^2/n}} \sim t_{n-1} \end{aligned}$$

In a test:

① Find $\tilde{\mu}$ & $\hat{\sigma}^2$

② Discrepancy measure:

$$t^* = \frac{\hat{\mu} - \mu_0}{\hat{\sigma}/\sqrt{n}} \sim t_{n-1}$$

$H_0: \mu = \mu_0$

③ p-value:

$$p\text{-val} = 2P(t_{n-1} > |t^*|) \rightarrow \text{2-sided } t\text{-test}$$

Single Treatment Confidence Interval

Pivotal quantity: $f(\text{data}, \text{unknown param from known distr.})$

When μ is param. of interest:

$$(1) \quad \frac{\hat{\mu} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \Rightarrow \text{If } \sigma \text{ known}$$

or

$$(2) \quad \frac{\hat{\mu} - \mu}{\hat{\sigma}/\sqrt{n}} \sim t_{n-1} \Rightarrow \sigma \text{ unknown}$$

From these quantities, we can derive $(1-\alpha) \cdot 100$ C.I.

$$\text{C.I.} = \hat{\mu} \pm c \cdot \text{s.e.}(\hat{\mu}) \rightarrow \text{s.e.}(\hat{\mu}) = \frac{\hat{\sigma}}{\sqrt{n}}$$

$$\hookrightarrow P(|t_{n-1}| \leq c) = 1 - \alpha$$

Interpretation:

✗ NOT 95% prob. that C.I. contain $\mu \Rightarrow$ it either does or doesn't

✓ One of many intervals that covers the mean μ 95% of time

✓ If you sample data many times, mean μ will lie in C.I. 95% of time

Hypothesis test for mean

① Establish H_0 & H_a

② Discrepancy measure:

$$\frac{\hat{\mu} - \mu_0}{\hat{\sigma}/\sqrt{n}} \sim t_{n-1}$$

③ Calculate p-value:

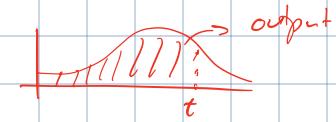
a) $H_a: \mu \neq 0 \Rightarrow p\text{-val} = 2 P(t_{n-1} > t_{\text{obs}})$

b) $H_a: \mu > 0 \Rightarrow p\text{-val} = P(t_{n-1} > t_{\text{obs}})$

c) $H_a: \mu < 0 \Rightarrow p\text{-val} = P(t_{n-1} < t_{\text{obs}})$

In R:

$\text{pt}(t, df)$ is CDF
of t-distr.



④ Statistical interpretation of p-val

Write in degree of evidence against H_0 . Close to zero \Rightarrow more evidence

⑤ Conclusion in context of study

Note that you cannot compare p-vals! Sample size is huge factor

C.I. will reflect results of hypothesis test. If $p\text{-val} < 0.05$, H_0 is 95% C.I.

COMPARATIVE EXPERIMENTAL PLANS

Two Treatments w/out Blocking

Model: $Y_{ij} = \mu + \tau_i + \beta_{ij}$, $i = 1, 2, j = 1 \dots n_i$, $\beta_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$

\uparrow Response \uparrow Overall mean \nearrow effect of treatment i \nearrow unit j 's error in treatment i

Model is subject to identifiability condition: $n_1 \tau_1 + n_2 \tau_2 = 0$

↳ Forces:

$$\begin{array}{c} \mu \\ \hline \left. \begin{array}{c} \tau_2 \\ \vdots \\ \tau_1 \\ \vdots \\ \mu \end{array} \right\} \end{array}$$

↳ If $n_1 = n_2 \Rightarrow$ balanced plan. O.L., unbalanced

↳ Importance:

A: Notation:

$$Y_{it} = \sum_{j=1}^{n_i} Y_{ij}, \quad \overline{Y}_{it} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}, \quad Y_{++} = \sum_{i=1}^2 \sum_{j=1}^{n_i} Y_{ij}, \quad \overline{Y}_{++} = \frac{1}{n_1 + n_2} Y_{++}$$

\uparrow Summing on other variable \checkmark mean

B: Estimability of parameters

We could write the following:

$$\begin{aligned} Y_{ij} &= \mu + \tau_i + R_{ij} = (\mu + \alpha) + (\tau_i - \alpha) + R_{ij} \\ &= \mu^* + \tau_i^* + R_{ij} \end{aligned}$$

This makes μ & τ_i hard to estimate. Condition prevents this.
b/c:

$$n_1 \tau_1^* + n_2 \tau_2^* = n_1 \tau_1 + n_2 \tau_2 - 2\alpha \neq 0$$

$\therefore \tau_i$ is strictly treatment effect

C: Interpretability:

We want $E[\bar{Y}_{++}] = \mu \Rightarrow$ pop. avg. of response variable.

$$\begin{aligned} E[\bar{Y}_{++}] &= \frac{1}{n_1 + n_2} \sum_{i=1}^2 \sum_{j=1}^{n_i} E[Y_{ij}] \\ &= \frac{1}{n_1 + n_2} \sum_{i=1}^2 \sum_{j=1}^{n_i} (\mu + \tau_i) \\ &= \frac{1}{n_1 + n_2} \sum_{i=1}^2 n_i (\mu + \tau_i) \\ &= \frac{1}{n_1 + n_2} (n_1 \mu + n_1 \tau_1 + n_2 \mu + n_2 \tau_2) \\ &= \frac{1}{n_1 + n_2} (n_1 + n_2) \mu + (n_1 \tau_1 + n_2 \tau_2) \\ &= \mu + (n_1 \tau_1 + n_2 \tau_2) \end{aligned}$$

If we want $E[\bar{Y}_{++}] = \mu \Rightarrow n_1 \tau_1 + n_2 \tau_2 = 0$

We are interested in diff. b/w treatments:

$$\theta = \tau_1 - \tau_2$$

which is:

$$\begin{aligned} E[\bar{Y}_{1+} - \bar{Y}_{2+}] &= E[\bar{Y}_{1+}] - E[\bar{Y}_{2+}] \\ &= \mu + \tau_1 - \mu - \tau_2 \\ &= \tau_1 - \tau_2 \end{aligned}$$

Thus $\tilde{\theta} = \bar{Y}_{1+} - \bar{Y}_{2+}$ is unbiased estimator of θ

Next goal: estimate μ & τ_i

↳ Proof: We know $Y_{ij} = \mu + \tau_i + R_{ij} \Rightarrow R_{ij} = Y_{ij} - \mu - \tau_i$

Loss function:

$$\min_{\mu, \tau_1, \tau_2} \sum_{i=1}^2 \sum_{j=1}^{n_i} R_{ij}^2, \text{ subject to } n_1 \tau_1 + n_2 \tau_2 = 0$$

We use method of Lagrange multipliers to solve this:

$$\min_{\mu, \tau_1, \tau_2, \lambda} \left[L(\mu, \tau_1, \tau_2) + \lambda (n_1 \tau_1 + n_2 \tau_2) \right]$$

$$\Rightarrow \min_{\mu, \tau_1, \tau_2, \lambda} \left[\sum_{i=1}^2 \sum_{j=1}^{n_i} (y_{ij} - \mu - \hat{\tau}_i)^2 + \lambda (n_1 \tau_1 + n_2 \tau_2) \right]$$

$\underbrace{\quad}_{Q}$

Thus:

$$① \frac{\partial Q}{\partial \mu} = -2 \sum_{i=1}^2 \sum_{j=1}^{n_i} (y_{ij} - \hat{\mu} - \hat{\tau}_i) = 0$$

$$② \frac{\partial Q}{\partial \tau_i} = \sum_{j=1}^{n_i} (-2(y_{ij} - \hat{\mu} - \hat{\tau}_i) + n_i \hat{\lambda}) = 0$$

$$③ \frac{\partial Q}{\partial \lambda} = n_1 \hat{\tau}_1 + n_2 \hat{\tau}_2 = 0$$

Solving ①

$$y_{++} - (n_1 + n_2) \hat{\mu} - n_1 \hat{\tau}_1 - n_2 \hat{\tau}_2 = 0$$

$$y_{++} - (n_1 + n_2) \hat{\mu} = 0$$

$$\hat{\mu} = \frac{y_{++}}{n_1 + n_2} = \bar{y}_{++}$$

o b/c ③

Solving ②

$$-2(y_{1+} - n_1 \hat{\mu} - n_1 \hat{\tau}_1) + n_1 \hat{\lambda} = 0$$

$$\Rightarrow -2 \sum_{i=1}^2 (y_{i+} - n_i \hat{\mu} - n_i \hat{\tau}_i) + \sum_{i=1}^2 n_i \hat{\lambda} = 0$$

) Common trick: add treatment effect

$$\Rightarrow y_{1+} - n_1 \hat{\mu} - n_1 \hat{\tau}_1 + y_{2+} - n_2 \hat{\mu} - n_2 \hat{\tau}_2 + \frac{1}{2} n_1 \hat{\lambda} + \frac{1}{2} n_2 \hat{\lambda} = 0$$

$$\Rightarrow y_{1+} + y_{2+} - (n_1 + n_2) \hat{\mu} + \frac{1}{2} \hat{\lambda} (n_1 + n_2) = 0$$

$$\Rightarrow \underbrace{y_{++}}_{y_{++}} - y_{++} + \frac{1}{2} \hat{\lambda} (n_1 + n_2) = 0$$

$$\Rightarrow \frac{1}{2} \hat{\lambda} (n_1 + n_2) = 0$$

$$\therefore \hat{\lambda} = 0$$

Reducing ②

$$\begin{aligned} \sum_{j=1}^{n_i} (-2(y_{ij} - \hat{\mu} - \hat{\tau}_i)) &= 0 \\ \Rightarrow y_{i+} - n_i \hat{\mu} - n_i \hat{\tau}_i &= 0 \\ \Rightarrow \hat{\tau}_i &= \hat{\mu} - \bar{y}_{i+} = \bar{y}_{i+} - \bar{y}_{i+} \end{aligned}$$

Variance estimation:

① Variance w/in treatment i :

$$\hat{\sigma}_i^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (y_{ij} - \hat{\mu} - \hat{\tau}_i)^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i+})^2$$

② Pooled/overall variance

$$\hat{\sigma}^2 = \frac{1}{n_1 + n_2 - 2} \sum_{i=1}^2 \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i+})^2 = \frac{(n_1 - 1) \hat{\sigma}_1^2 + (n_2 - 1) \hat{\sigma}_2^2}{n_1 + n_2 - 2}$$

Statistical inference

Note that $\bar{Y}_{i+} \sim N(\mu + \tau_i, \sigma^2/n_i)$

\hookrightarrow Proof:

① Normality:

$$\bar{Y}_{i+} = \frac{1}{n_i} (y_{i1} + \dots + y_{in_i}) \Rightarrow \text{indep. L.C. of normal R.V.} \quad \therefore \bar{Y}_{i+} \text{ is normal}$$

② Mean:

$$E[\bar{Y}_{i+}] = \frac{1}{n_i} \sum_{j=1}^{n_i} E[\mu + \tau_i + \varepsilon_{ij}] = \mu + \tau_i$$

③ Variance:

$$\text{Var}[\bar{Y}_{i+}] = \frac{1}{n_i^2} \sum_{j=1}^{n_i} \text{Var}[\mu + \tau_i + \varepsilon_{ij}] = \frac{\hat{\sigma}_i^2}{n_i}$$

Therefore:

$$\bar{Y}_{1+} - \bar{Y}_{2+} \sim N(\tau_1 - \tau_2, \sigma^2(\frac{1}{n_1} + \frac{1}{n_2}))$$

We also know that

$$\frac{(n_1 + n_2 - 2) \tilde{\sigma}^2}{\sigma^2} \sim \chi^2_{n_1 + n_2 - 2}$$

Some argument as
 $\frac{\bar{\mu} - \mu}{\sigma/\sqrt{n}} \sim t_{n-1}$ derivation

Thus:

$$\frac{(\bar{Y}_{1+} - \bar{Y}_{2+}) - \theta}{\tilde{\sigma} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1 + n_2 - 2}$$

Pooled
std.
dev.

Use this to make C.I. & tests

This all relies on $\beta_{ij} \sim N(0, \sigma^2)$

Two Treatments w/ Blocking

Complete randomized block design: all treatments randomly assigned within each block, b blocks

Model:

$$Y_{ij} = \mu + \tau_i + \beta_j + \epsilon_{ij}, \quad \epsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2), \quad i=1, 2, \quad j=1, \dots, b$$

↑ block effect

subject to: $\tau_1 + \tau_2 = 0, \quad \sum_{j=1}^b \beta_j = 0 \quad \left. \right\} \text{identifiability cond.}$

Parameter estimation:

$$\text{Obj.} \quad \min_{\mu, \tau, \beta} \sum_{i=1}^2 \sum_{j=1}^b (y_{ij} - \mu - \tau_i - \beta_j)^2 \quad \text{subj.} \quad \begin{cases} \tau_1 + \tau_2 = 0 \\ \sum_{j=1}^b \beta_j = 0 \end{cases}$$

Use the same methods (Lagrange):

$$\min_{\mu, \tau, \beta, \lambda_1, \lambda_2} \left[\sum_{i=1}^2 \sum_{j=1}^b (y_{ij} - \mu - \tau_i - \beta_j)^2 + \lambda_1 (\tau_1 + \tau_2) + \lambda_2 \sum_{j=1}^b \beta_j \right]$$

Set derivatives:

$$\frac{\partial Q}{\partial \mu} = 0, \quad \frac{\partial Q}{\partial \tau_i} = 0, \quad \frac{\partial Q}{\partial \beta_j} = 0, \quad \frac{\partial Q}{\partial \lambda_1} = 0$$

Estimates:

$$\hat{\mu} = \bar{y}_{++}, \quad \hat{\tau}_i = \bar{y}_{i+} - \bar{y}_{++}, \quad \hat{\beta}_j = \bar{y}_{+j} - \bar{y}_{++}$$

Estimating variance:

$$\begin{aligned}
 \hat{\sigma}^2 &= \frac{1}{b-1} \sum_{i=1}^2 \sum_{j=1}^b (y_{ij} - \hat{\mu} - \hat{\tau}_i - \hat{\beta}_j)^2 \\
 &= \frac{1}{b-1} \sum_{i=1}^2 \sum_{j=1}^b (y_{ij} - \bar{y}_{i+} - \bar{y}_{+j} + \bar{y}_{++})^2 \\
 &= \frac{1}{2} \frac{\sum_{j=1}^b (\bar{d}_j - \bar{d})^2}{b-1}
 \end{aligned}$$

$\bar{d}_j = y_{ij} - y_{2j}$
 $\bar{d} = \frac{1}{b} \sum_{j=1}^b \bar{d}_j$
 * Prove this *

Why do we have $b-1$ d.f.?

↳ D.F. = sample size - # of estimated params. + # of constraints

$\underbrace{2 \text{ perf}}_{\text{per block}} = 2b - (1 + 2 + b) + 2$

In ex. = $b-1$
 = $n_1 + n_2 - 2$

Statistical inference

$\theta = \tau_1 - \tau_2 \Rightarrow$ param. of interest. Define distn. of $\hat{\theta}$

Claim: $\hat{\theta} \sim N(\tau_1 - \tau_2, \frac{2\sigma^2}{b})$

Proof: ① Normality

$$\hat{\theta} = \bar{Y}_{1+} - \bar{Y}_{2+} = \frac{1}{b} \sum_{j=1}^b Y_{1j} - \frac{1}{b} \sum_{j=1}^b Y_{2j} = \frac{1}{b} \sum_{j=1}^b (Y_{1j} - Y_{2j})$$

This is a linear combo of indep. normal R.V. $\Rightarrow \hat{\theta}$ is Gaussian

② Mean

$$\begin{aligned}
 E[\hat{\theta}] &= \frac{1}{b} \sum_{j=1}^b (E[Y_{1j}] - E[Y_{2j}]) \\
 &= \frac{1}{b} \sum_{j=1}^b (\tau_1 - \tau_2) \\
 &= \tau_1 - \tau_2
 \end{aligned}$$

③ Variance

$$\begin{aligned}
 \text{Var}[\hat{\theta}] &= \frac{1}{b^2} \sum_{j=1}^b (\text{Var}[Y_{1j}] + \text{Var}[Y_{2j}]) \\
 &= \frac{1}{b^2} \sum_{j=1}^b (2\sigma^2)
 \end{aligned}$$

$$\text{Var}(a-b) = \text{Var}(a) + \text{Var}(b) \\
 \downarrow \\
 + \text{Cov}(a, b)$$

$$= \frac{2\sigma^2}{b}$$

Let pivotal quantity be:

$$\frac{\bar{\theta} - (\tau_1 - \tau_2)}{\sigma \sqrt{2/b}} \sim N(0, 1)$$



indep.

We know

$$\frac{(b-1)\bar{\sigma}^2}{\sigma^2} \sim \chi^2_{b-1}$$

$$\Rightarrow \frac{\bar{\theta} - (\tau_1 - \tau_2)}{\bar{\sigma} \sqrt{2/b}} \sim t_{b-1}$$

Use this for C.I. & tests.

EXPERIMENTAL PLANS: ≥ 2 TREATMENTS

Introduction

Experimental design

Completely randomized

No blocking

Randomized block design

Each treatment replicates in each block, random assignment

Completely Randomized Design - Balanced

Balanced plan: same # of units receive each treatment

Goal: develop methodology to test $H_0: \mu_a = \mu_b = \mu_c = \dots = \mu_t$

① MODEL

Assuming balanced plan:

$$Y_{ij} = \mu + \tau_i + R_{ij}, \quad i=1, \dots, t, \quad j=1, \dots, r, \quad R_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$$

of treatments
↓
t
of repl.

Constraint: $\sum_{i=1}^t \tau_i = 0 \quad \leftarrow \text{exact same arg. as previous constraint}$

Meaning:

μ : mean response across all treatments

τ_i : effect of treatment i relative to μ

ϵ_{ij} : random error

Constraint leads us to claim $E[\bar{Y}_{++}] = \mu$

Proof:

$$\begin{aligned} E[\bar{Y}_{++}] &= \frac{1}{r \cdot t} \sum_{i=1}^t \sum_{j=1}^r E[Y_{ij}] \\ &= \frac{1}{r \cdot t} \sum_{i=1}^t \sum_{j=1}^r (\mu + \tau_i) \\ &= \frac{1}{r \cdot t} \sum_{i=1}^t (\mu + \tau_i) \\ &= \frac{t\mu}{t} + \frac{1}{t} \sum_{i=1}^t \tau_i \xrightarrow{0 \text{ b/c of constraint}} \\ &= \mu \end{aligned}$$

We care about treatment effect τ_i :

$$E[\bar{Y}_{it}] = \mu + \tau_i \Rightarrow \tau_i = \underbrace{E[\bar{Y}_{it}]}_{\text{Unbiased estimator of } \tau_i} - \overline{E[\bar{Y}_{++}]} \quad \text{Unbiased estimator of } \tau_i$$

② PARAMETER ESTIMATION

Loss func:

$$\sum_{i=1}^t \sum_{j=1}^r (y_{ij} - (\mu + \hat{\tau}_i))^2 + \lambda \sum_{i=1}^t \hat{\tau}_i$$

Solving this leads to:

$$\hat{\mu} = \bar{Y}_{++} \Rightarrow \hat{\mu} = \bar{Y}_{++}$$

$$\hat{\tau}_i = \bar{Y}_{it} - \bar{Y}_{++} \Rightarrow \hat{\tau}_i = \bar{Y}_{it} - \bar{Y}_{++}$$

$$\hat{\tau}_i - \hat{\tau}_j = \bar{Y}_{it} - \bar{Y}_{jt} \Rightarrow \hat{\theta} = \bar{Y}_{it} - \bar{Y}_{jt}$$

Variance estimator:

$$\hat{\sigma}^2 = \frac{1}{rt - (1+t) + 1} \sum_{i,j} (y_{ij} - \hat{\mu} - \hat{\tau}_i)^2 = \frac{1}{t(r-1)} \sum_{i,j} (y_{ij} - \hat{\mu} - \hat{\tau}_i)^2$$

Sample size # of params * of constraints

Note:

$$\hat{\sigma}^2 = \frac{(\nu-1) \hat{\sigma}_1^2 + \dots + (\nu-1) \hat{\sigma}_t^2}{\nu-1}$$

③ DISTRIBUTION

We know $\bar{Y}_{it} \stackrel{\text{ind.}}{\sim} N(\mu + \tau_i, \sigma^2/r)$

$$\Rightarrow \tilde{\theta} = \tilde{\tau}_i - \tilde{\tau}_j = \bar{Y}_{it} - \bar{Y}_{jt} \sim N(\tau_i - \tau_j, \frac{2\sigma^2}{r}) \Rightarrow \text{CONTRAST}$$

$$\Rightarrow \frac{\tilde{\theta} - (\tau_i - \tau_j)}{\sigma \sqrt{2/r}} \sim N(0, 1)$$

$$\text{We also know } U = \frac{\nu(\nu-1) \tilde{\sigma}^2}{\sigma^2} \sim \chi^2_{\nu(\nu-1)}$$

$$\Rightarrow \frac{\tilde{\theta} - (\tau_i - \tau_j)}{\tilde{\sigma} \sqrt{2/r}} \sim t_{\nu(\nu-1)}$$

Completely Randomized Design: Generalizing via Contrasts

Contrast defn: linear combo. of treatment effects s.t. weights sum to 0

Mathematically: $a_1, \dots, a_t \in \mathbb{R}$ (weights), $\tau_1, \dots, \tau_t \in \mathbb{R}$ (treatment effects)

$$\textcircled{1} \quad \sum_i a_i \tau_i \quad (\text{linear combo})$$

$$\textcircled{2} \quad \sum_{i=1}^t a_i = 0$$

Ex:// Identify contrast

$$1) \quad \tau_1 - \frac{\tau_2 + \tau_3}{2} \Rightarrow \text{Yes! L.C., } \sum a_i = 0$$

$$2) \quad \tau_i^2 - \tau_j \Rightarrow \text{No, not L.C.}$$

$$3) \quad \frac{\tau_1 + \tau_2}{2} - \frac{\tau_3 + \tau_4 + \tau_5}{3} \Rightarrow \text{Yes, } \sum a_i = 0$$

Ex:// Write contrasts to compare:

1) Treatment A vs. B, C, D

Trick: if A vs. B.

$$\bar{\tau}_A = \frac{\bar{\tau}_B + \bar{\tau}_C + \bar{\tau}_D}{3} \Rightarrow \begin{array}{l} \text{1. Make coef. of A\&B sum to 1!} \\ \text{2. } A-B \end{array}$$

2) A & B vs. C & D

$$\frac{\bar{\tau}_A + \bar{\tau}_B}{2} - \frac{\bar{\tau}_C + \bar{\tau}_D}{2} \Rightarrow (\bar{\tau}_A + \bar{\tau}_B) - (\bar{\tau}_C + \bar{\tau}_D)$$

Find the estimator of a generic contrast:

1. Parameter: $\theta = \sum_{i=1}^t a_i \tau_i, \sum_{i=1}^t a_i = 0$

2. Estimate:

$$\begin{aligned} \hat{\theta} &= \sum_{i=1}^t a_i \hat{\tau}_i = \sum_{i=1}^t a_i (\bar{y}_{i+} - \bar{y}_{++}) \\ &= \sum_{i=1}^t a_i \bar{y}_{i+} - \bar{y}_{++} \sum_{i=1}^t a_i \xrightarrow{=} 0 \\ &= \sum_{i=1}^t a_i \bar{y}_{i+} \end{aligned}$$

3. Estimator:

$$\tilde{\theta} = \sum_{i=1}^t a_i \bar{Y}_{i+}$$

STDEV, not var.

What is distn. of $\tilde{\theta}$? Claim: $\tilde{\theta} \sim G(\theta, \sigma \sqrt{\sum_{i=1}^t a_i^2 / r})$

Proof:

① Normality: $\tilde{\theta}$ is sum of indep. normal R.V.

② Expectation:

$$\begin{aligned} E[\tilde{\theta}] &= \sum_{i=1}^t a_i E[\bar{Y}_{i+}] = \sum_{i=1}^t a_i (\mu + \tau_i) \\ &= \mu \sum_{i=1}^t a_i \xrightarrow{=} 0 + \sum_{i=1}^t a_i \tau_i \\ &= \theta \end{aligned}$$

③ Variance:

$$\text{Var}[\tilde{\theta}] = \sum_{i=1}^t a_i^2 \text{Var}(\bar{Y}_{i+}) = \sum_{i=1}^t a_i^2 \frac{\sigma^2}{r}$$

Discrepancy measure for generic contrast:

$$\frac{\hat{\theta} - \theta}{\tilde{\sigma} \sqrt{\sum a_i^2 / r}} \sim t_{t(r-1)}$$

We can test diff. b/w generic combos of treatments

Completely Randomized Design: ANOVA

Goal: develop methodology to test $H_0: \mu_a = \mu_b = \mu_c = \dots = \mu_r$

Idea: Construct:

A: $\tilde{\sigma}_r^2$: estimator for σ^2 , unbiased regardless if H_0 true/not

From model, we know:

$$\tilde{\sigma}_r^2 = \frac{1}{t(r-1)} \sum_{i,j} (y_{ij} - \bar{y}_{it})^2 \Rightarrow \text{unbiased}$$

B: $\tilde{\sigma}_t^2$: II, unbiased $\Leftrightarrow H_0$ is true

From model: $\bar{y}_{it} \sim N(\mu + \tau_i, \sigma^2/r)$

If H_0 is true, then $\tau_1 = \dots = \tau_t = 0$.

$$\Rightarrow \bar{y}_{it} \sim N(\mu, \sigma^2/r)$$

Estimate variance:

$$\hat{\sigma}_r^2 = \frac{\sum_t (\bar{y}_{it} - \bar{\bar{y}}_{tr})^2}{t-1} = \frac{1}{t-1} \sum \hat{\tau}_i^2$$

$$\therefore \hat{\sigma}_{tr}^2 = \frac{r}{t-1} \sum \hat{\tau}_i^2 \Rightarrow \tilde{\sigma}_{tr}^2 = \frac{r}{t-1} \sum \tilde{\tau}_i^2$$

treatment mean square

\therefore Following estimator:

$$F = \frac{\tilde{\sigma}_t^2}{\tilde{\sigma}_r^2} \Rightarrow \begin{array}{ll} \text{If } H_0 \text{ is true} \Rightarrow \text{value } \approx 1 & \\ \text{If } H_0 \text{ is false} \Rightarrow B \uparrow \rightarrow F \uparrow & \end{array}$$

↑ residual mean square

Distribution of F:

We know:

$$\textcircled{1} \quad \frac{t(-1) \tilde{\sigma}_r^2}{\sigma^2} \sim \chi^2_{r(t-1)}$$

$$\textcircled{2} \quad \frac{(t-1) \tilde{\sigma}_t^2}{\sigma^2} \sim \chi^2_{t-1}$$

Also:

$$X \sim \chi^2_{df_X}, \quad Y \sim \chi^2_{df_Y} \Rightarrow F = \frac{X/df_X}{Y/df_Y} \sim F_{df_X, df_Y}$$

Thus:

$$\frac{\tilde{\sigma}_t^2}{\tilde{\sigma}_r^2} \sim F_{t-1, r(t-1)}$$

ANOVA & Sum of Squares

Theorem: sum of squared errors can be decomposed like so:

$$\sum_{i,j} (y_{ij} - \bar{y}_{++})^2 = \underbrace{\sum_{i,j} (y_{ij} - \bar{y}_{it})^2}_{\text{Total sum of squares}} + r \underbrace{\sum_i (\bar{y}_{it} - \bar{y}_{++})^2}_{\text{Residual/error sum of squares}} + \underbrace{r \sum_i (\bar{y}_{it} - \bar{y}_{++})^2}_{\text{Treatment sum of squares}}$$

Proof:

$$\begin{aligned} \sum_{i=1}^t \sum_{j=1}^r (y_{ij} - \bar{y}_{++})^2 &= \sum_{i=1}^t \sum_{j=1}^r (y_{ij} - \bar{y}_{it} + \bar{y}_{it} - \bar{y}_{++})^2 \\ &= \sum_{i=1}^t \sum_{j=1}^r (y_{ij} - \bar{y}_{it})^2 + \sum_{i=1}^t \sum_{j=1}^r (\bar{y}_{it} - \bar{y}_{++})^2 - 2 \sum_{i,j} (\bar{y}_{it} - \bar{y}_{++})(y_{ij} - \bar{y}_{it}) \\ &= \sum_{i=1}^t \sum_{j=1}^r (y_{ij} - \bar{y}_{it})^2 + \sum_{i=1}^t \sum_{j=1}^r (\bar{y}_{it} - \bar{y}_{++})^2 \quad \square \end{aligned}$$

$$\begin{aligned} \sum_i (\bar{y}_{it} - \bar{y}_{++}) &= -t \bar{y}_{++} + \sum_{i=1}^t \bar{y}_{it} \\ &= -t \bar{y}_{++} + t \bar{y}_{++} \\ &= 0 \end{aligned}$$

Intuition: if the treatment S.S. \gg error S.S. \Rightarrow treatments are diff.

- SST (total): total variation in response
- SSTR (treatment): variability b/w treatment groups
- SSE (error): variability w/in group. If large \rightarrow harder to detect diff. b/w groups

ANOVA table:

	Sum of squares	DF	Mean square	F_{obs}
Treatment	$r \sum_i (\bar{y}_{it} - \bar{y}_{++})^2$	$t-1$	$= \frac{SSTr}{df_{tr}}$	$\frac{MS_{tr}}{MS_r}$
Residual	$\sum_{i,j} (y_{ij} - \bar{y}_{i+})^2$	$t(r-1)$	$= \frac{SSR}{df_r}$	
Total	$\sum_{i,j} (y_{ij} - \bar{y}_{++})^2$	$tr-1$		<p>Be able to construct w/ only summary stats</p>

Intuition: large diff. in treatment \Rightarrow large F_{obs} ($MS_{tr} > MS_r$)

ANOVA test:

$$\left\{ \begin{array}{l} H_0: \tau_1 = \dots = \tau_t = 0/\mu \\ H_a: \text{not the same} \end{array} \right. \Rightarrow \begin{array}{l} \textcircled{1} \text{ Calculate } F_{\text{obs}} \\ \textcircled{2} \text{ p-value: } P(F_{t-1, t(r-1)} \geq F_{\text{obs}}) \end{array}$$

Note: If $t=2$: $F_{\text{obs}} = (F_{\text{obs}})^2 \Rightarrow$ know proof

After this, perform post-hoc analysis to see which treatments are actually diff.

Completely Randomized Design - Unbalanced

Defn: t treatments assigned to group randomly but # of replications / treatment diff

① MODEL

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

$$\epsilon_{ij} \sim N(0, \sigma^2), \quad \sum_{i=1}^t \tau_i = 0$$

Interpretation:

μ : weighted avg. of treatment effects

$\mu + \tau_i$: average resp. for treatment i

② PARAMETER ESTIMATION

$$\min_{\mu, \tau_1, \dots, \tau_t, \lambda} \left[\sum_{i=1}^t \sum_{j=1}^{r_i} (y_{ij} - \mu - \tau_i)^2 + \lambda \sum_{i=1}^t r_i \tau_i \right]$$

$$\hat{\mu} = \bar{y}_{++}, \quad \hat{\tau}_i = \bar{y}_{ir} - \bar{y}_{++}$$

$$\hat{\mu} = \bar{Y}_{++}, \quad \hat{\tau}_i = \bar{Y}_{ir} - \bar{Y}_{++}$$

Main change:

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^t \sum_{j=1}^{r_i} (y_{ij} - \bar{y}_{i+})^2}{\sum_{i=1}^t (r_i - 1)} \xrightarrow{\text{Equiv. to } y_{ij} - \bar{y}_{i+}}$$

d.f.

Why is $\sum_{i=1}^t (r_i - 1)$ d.f.?

$$\sum_{i=1}^t (r_i - 1) = \sum_{i=1}^t r_i - (t + 1) + 1$$

of obs. # of param. constraint

③ ANOVA

	Sum of squares	DF	Mean square	F _{obs}
Treatment	$\sum_{i,j} (\bar{y}_{it} - \bar{y}_{++})^2$	$t-1$	SST_r / df_r	$\frac{MS_{tr}}{MS_r}$
Residual	$\sum_{i,j} (y_{ij} - \bar{y}_{i+})^2$	$\sum_{i=1}^t (r_i - 1)$	SSR / df_r	
Total	$\sum_{i,j} (y_{ij} - \bar{y}_{++})^2$	$\sum_{i=1}^t r_i - 1$		Equiv. if r_i is same $\forall i$ to par. ANOVA

Ex // Fill in ANOVA for following data:

	T ₁	T ₂	T ₃	Overall
Avg.	3.64	3.40	3.21	3.41
Variance	0.123	0.067	0.105	0.125
n	13	14	14	41

① Estimate params:

$$\hat{\mu} = \bar{y}_{\text{fit}} = 3.91, \quad \hat{\tau}_1 = \bar{y}_{1+} - \bar{y}_{7+} = 0.23, \quad \hat{\tau}_2 = -0.01, \quad \hat{\tau}_3 = -0.2$$

② Sum of squares:

$$\begin{aligned}\text{Total sum of squares} &= \text{d.f.}_T \times \text{overall variance} \\ &= (13 + 14 + 14 - 1) \times 0.125 \\ &= 5\end{aligned}$$

$$SS_{\text{tr}} = \sum_{i=1}^t \sum_{j=1}^{r_i} \left(\underbrace{\bar{y}_{ij}}_{\hat{\tau}_i} - \bar{y}_{\text{fit}} \right)^2 = \sum_{i=1}^t r_i \hat{\tau}_i^2 = 1.249$$

$$SSE = SST - SS_{\text{tr}} = 3.751$$

③ D.f.

$$\text{d.f.}_T = 2, \quad \text{d.f.}_{\text{res}} = \text{d.f.}_T - \text{d.f.}_{\text{tr}} = 40 - 2 = 38$$

④ Fill in MS & Fobs.

How to do tests for contrasts:

① Estimate

$$\hat{\theta} = \sum a_i \hat{\tau}_i = \sum a_i \bar{y}_{ij}$$

② Distr.:

$$\begin{aligned}\bar{y}_{ij} &\sim N(\mu + \tau_i, \frac{\sigma^2}{r_i}) \\ \Rightarrow \hat{\theta} = \sum a_i \bar{y}_{ij} &\sim N\left(\sum_{i=1}^t a_i \tau_i, \sigma^2 \sum_{i=1}^t \frac{a_i^2}{r_i}\right) \\ \therefore \hat{\theta} \pm t^* \hat{\sigma} \sqrt{\sum_{i=1}^t \frac{a_i^2}{r_i}} &\text{ is 95% CI.}\end{aligned}$$

Randomized Block Design

Defn: multiple treatments assigned randomly w/in blocks

① MODEL

$$Y_{ij} = \mu + \tau_i + \beta_j + \beta_{ij}, \quad i=1, \dots, t, \quad j=1, \dots, b$$

$$\beta_{ij} \sim N(0, \sigma^2), \quad \sum_{i=1}^t \tau_i = \sum_{j=1}^b \beta_j = 0$$

Interpretations:

μ : overall avg. across treat. & blocks.

τ_i : i^{th} treatment effect

β_j : j^{th} block effect

Additive model: no interaction b/w block & treatmt

① Treatment effect does not affect block (τ_i , not τ_{ij})

② Diff. b/w treatments is same \forall blocks (β_j , not β_{ij})

③ Residual variance constant across treat. & blocks

② PARAM. ESTIMATES

$$\min_{\mu, \tau_1, \dots, \tau_t, \beta_1, \dots, \beta_b} \left[\sum_{i=1}^t \sum_{j=1}^b (y_{ij} - \mu - \tau_i - \beta_j)^2 + \lambda_1 \sum_i \tau_i + \lambda_2 \sum_j \beta_j \right]$$

Estimate:

$$\hat{\mu} = \bar{y}_{++}, \quad \hat{\tau}_i = \bar{y}_{i+} - \bar{y}_{++}, \quad \hat{\beta}_j = \bar{y}_{+j} - \bar{y}_{++}$$

Variance estimate:

$$\begin{aligned} \textcircled{i} \text{ D.f.} &: \text{total #} - \# \text{ of parsns} + \# \text{ of const.} \\ &= tb - (t+1+b) + 2 \\ &= (t-1)(b-1) \end{aligned}$$

③ TEST

Goal: test treatment effects simultaneously

Sum of squares decomp:

$$\hat{\sigma}^2 = \frac{1}{(t-1)(b-1)} \sum_{i,j} (y_{ij} - \bar{y}_{i+} - \bar{y}_{+j} + \bar{y}_{++})^2$$

$$\sum_{i,j} (y_{ij} - \bar{y}_{++})^2 = \sum_{i,j} (y_{ij} - \bar{y}_{it} - \bar{y}_{rj} + \bar{y}_{++})^2 + b \sum_{i=1}^t (\bar{y}_{it} - \bar{y}_{r+})^2 + t \sum_{j=1}^b (\bar{y}_{rj} - \bar{y}_{++})^2$$

SS_{total} SS_{error} SS_{treatment} SS_{block}
 pulled out of SSE

Biased variance estimate if $H_0: \tau_i = \dots = \tau_t = 0$

Assume H_0 true $\Rightarrow \bar{Y}_{it} \sim N(\mu, \sigma^2/b)$

$$\therefore \hat{\sigma}_t^2 = \frac{b}{t-1} \sum_{i=1}^t (\bar{Y}_{it} - \bar{Y}_{++})^2$$

$$\therefore \frac{SS_{treatment}}{\hat{\sigma}^2} \sim \chi^2_{t-1}$$

$$F_{obs} = \frac{MS_T}{MS_r} \sim F_{t-1, (t-1)(b-1)}$$

ANOVA table:

	Sum of squares	DF	Mean square	F _{obs}
Treatment	$b \sum_i (\bar{y}_{it} - \bar{y}_{++})^2$	$t-1$	$SSTr/df_{tr}$	$\frac{MS_{tr}}{MS_r}$
Blocks	$t \sum_j (\bar{y}_{rj} - \bar{y}_{++})^2$	$b-1$	SS_B/df_B	
Residual	$\sum_{i,j} (y_{ij} - \bar{y}_{it} - \bar{y}_{rj} + \bar{y}_{++})^2$	$(t-1)(b-1)$	SS_R/df_r	
Total	$\sum_{i,j} (y_{ij} - \bar{y}_{++})^2$	$t(b-1)$		Compute via summary stats.

If blocking is important:

$$MSR \leftarrow \frac{SS_B + SS_{resid}}{df_B + df_{resid}}$$

Meaningful amt. of variance explained by blocking

Declining w/ contrasts:

$$\hat{\theta} = \sum_i a_i \bar{y}_{it} \sim N(\theta, \sigma^2 \frac{\sum a_i^2}{b})$$

$$U = \frac{(t-1)(b-1)\sigma^2}{\sigma^2} \sim \chi^2_{(t-1)(b-1)}$$

$$\Rightarrow \frac{\hat{\theta} - \theta}{\sigma_r \sqrt{\frac{\sum a_i^2}{b}}} \sim t_{(t-1)(b-1)} \Rightarrow U_{rc} \text{ for tests}$$

Why is blocking useful? Reducing residual error \rightarrow easier to detect treatment diff if effect exists.

↳ Blocking prevents treatments to be assigned to same type of units

FACTORIAL TREATMENT STRUCTURE & INTERACTION

Blocking acted like a nuisance variable: we don't care about it but it helps identify treatment diff.

Factorial treatment structure: compare combo of factor levels, each combo considered as a treatment.

Diff. b/w factorial vs. blocked design:

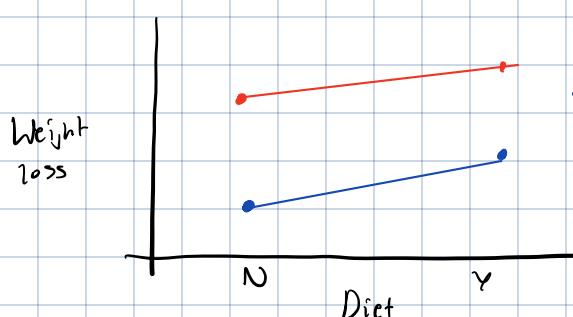
- ① Intent: factorial covers all factors, blocked only covers about treatment effects
- ② Randomization: blocked could only randomize treatments in a block, factorial randomizes everything.

Interaction

Defn: effect of 1 factor \leftarrow level of another factor

Interaction plot: all treatments in 1 plot

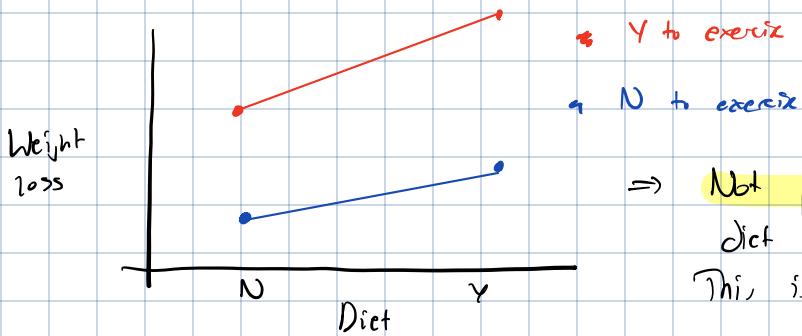
Ex://



* Y to exerciz
+ N to exerciz

\Rightarrow Parallel, so exercise doesn't impact weight loss regardless of diet

However



\Rightarrow Not parallel! If exercising, diet impact is diff.
This is an interaction

Two factor interaction plot: effect on response is same direction, but magnitude \uparrow

- i.e. \uparrow : \uparrow differs depending on other factor.
↳ crosses lines in interaction plot.

We can test interaction like a contrast if only 2 factors, w/ 2 levels:

$$\theta = (\bar{Y}_a - \bar{Y}_b) - (\bar{Y}_c - \bar{Y}_d)$$

↴ ↴
 Both treatments Both treatments
 have same factor have same factor
 level for factor X level for factor X

\Rightarrow Test of change in factor
X impacts response of
change factor Y

$$J = \frac{\hat{\theta} - \theta_0}{\hat{\sigma} \sqrt{\frac{1}{r} \sum a_{ii}^2}}$$

↳ Assuming balanced

Interaction Model

Model: $Y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ijk}, \epsilon_{ijk} \sim N(0, \sigma^2)$

Interpretations:

① Y_{ijk} : response of replicate k for factor A at level i, factor B at level j

② μ : mean response across all treatments

③ α_i : treatment effect for level i for factor A

④ β_j : \uparrow j \uparrow B

⑤ γ_{ij} : interaction effect b/w level i for factor A & level j for factor B

$$\text{Note: } E[Y_{ijk}] = \mu + \alpha_i + \lambda_j + \gamma_{ij} = \mu + \tau$$

$$\text{Num of params} = 1 + t_a + t_b + t_a t_b$$

Identifiability conditions:

$$\textcircled{1} \quad \sum_i \alpha_i = 0$$

$$\textcircled{2} \quad \sum_j \lambda_j = 0$$

$$\textcircled{3} \quad \sum_i \gamma_{ij} = 0 \quad \forall j \in [1, t_b] \quad \rightarrow \quad t_a + t_b - 1 \text{ constraints for } \gamma_{ij}$$

$$\textcircled{4} \quad \sum_j \gamma_{ij} = 0 \quad \forall i \in [1, t_a] \quad \rightarrow$$

Parameter estimation

\textcircled{1} Notation:

$$\bar{Y}_{+++} = \frac{1}{r t_a t_b} \sum_i \sum_j \sum_k Y_{ijk}, \quad \bar{Y}_{+jk} = \frac{1}{t_a t_b} \sum_i \sum_j Y_{ijk} \quad \begin{matrix} \text{avg. of } k^{\text{th}} \\ \text{replication} \\ \text{of treatments} \end{matrix}$$

\textcircled{2} Estimates:

$$\min_{\mu, \alpha_i, \lambda_j, \gamma_{ij}} \sum_i \sum_j \sum_k (y_{ijk} - \mu - \alpha_i - \lambda_j - \gamma_{ij})^2$$

yields:

$$\textcircled{1} \quad \hat{\mu} = \bar{Y}_{+++}$$

$$\textcircled{2} \quad \hat{\alpha}_i = \bar{Y}_{i++} - \bar{Y}_{+++}$$

$$\textcircled{3} \quad \hat{\lambda}_j = \bar{Y}_{+j+} - \bar{Y}_{+++}$$

$$\textcircled{4} \quad \hat{\gamma}_{ij} = \bar{Y}_{ij+} - \bar{Y}_{i++} - \bar{Y}_{+j+} + \bar{Y}_{+++} \Rightarrow \text{hard to interpret}$$

Also:

$$\hat{\sigma}^2 = \frac{1}{t_a t_b (r-1)} \sum_i \sum_j \sum_k (y_{ijk} - \hat{\mu} - \hat{\alpha}_i - \hat{\lambda}_j - \hat{\gamma}_{ij})^2$$

Blocking in interaction model

of replication of each treatment = # of blocks \Rightarrow 1 treatment per block

$$\therefore Y_{ijk} = \mu + \alpha_i + \gamma_j + \gamma_{ij} + \beta_k = \sum_{\substack{i \in [1, t_a] \\ j \in [1, t_b]}} (\bar{y}_{ij+} - \bar{y}_{+++})^2 \quad k \in [1, b]$$

Parameter estimation is same, w/ $\hat{\beta}_k = \bar{y}_{++k} - \bar{y}_{+++}$

Decomposition of treatment S.S.

Recall that:

$$SS_{\text{treat}} = b \sum_i \sum_j (\bar{y}_{ij+} - \bar{y}_{+++})^2$$

We want:

$$SS_{\text{treat}} = SS_A + SS_B + SS_{\text{int}}$$

↓ ↓ ↓
 factor S.S. interaction

Proof sketch:

$$\begin{aligned}
 b \sum_i \sum_j (\bar{y}_{ij+} - \bar{y}_{+++})^2 &= b t_b \sum_i (\bar{y}_{i++} - \bar{y}_{+++})^2 \quad] \rightarrow SS_A \\
 &\quad + b t_a \sum_j (\bar{y}_{+j+} - \bar{y}_{+++})^2 \quad] \rightarrow SS_B \\
 &\quad + b \sum_i \sum_j (\bar{y}_{ij+} - \bar{y}_{i++} - \bar{y}_{+j+} + \bar{y}_{+++})^2 \quad] \rightarrow SS_{\text{int}}
 \end{aligned}$$

↙ $\pm \bar{y}_{i++}, \bar{y}_{+j+}, \bar{y}_{++k} \text{ & } (a+b+c)^2$ identity

* Be able to show this *

$$\begin{aligned}
 \text{D.F. : } df_A + df_B + df_{\text{int}} &= (t_a - 1) + (t_b - 1) + (t_a - 1)(t_b - 1) \\
 &= t_a t_b - 1 \\
 &\hookrightarrow SS_{\text{treat}} \text{ D.F.}
 \end{aligned}$$

ANOVA table:

Source	Sum of squares	DF	Mean square
Treatments	$b \sum_{i,j} (\bar{y}_{ij+} - \bar{y}_{+++})^2$	$t_a t_b - 1$	$\frac{SS_T}{(t_a t_b - 1)}$
- Factor A	$t_b \sum_i (\bar{y}_{i++} - \bar{y}_{+++})^2$	$t_a - 1$	$\frac{SS_A}{(t_a - 1)}$
- Factor B	$t_a \sum_j (\bar{y}_{+j+} - \bar{y}_{+++})^2$	$t_b - 1$	$\frac{SS_B}{(t_b - 1)}$
- Interactions	$b \sum_{i,j} (\bar{y}_{ij+} - \bar{y}_{i++} - \bar{y}_{+j+} + \bar{y}_{+++})^2$	$(t_a - 1)(t_b - 1)$	$\frac{SS_I}{(t_a - 1)(t_b - 1)}$
Blocks	$t_a t_b \sum_k (\bar{y}_{++k} - \bar{y}_{+++})^2$	$b - 1$	$\frac{SS_B}{(b - 1)}$
Residual	$\sum_{i,j,k} (y_{ijk} - \bar{y}_{ij+} - \bar{y}_{+jk} + \bar{y}_{+++})^2$	$(t_a t_b - 1)(b - 1)$	$\frac{SS_R}{(t_a t_b - 1)(b - 1)}$
Total	$\sum_{i,j,k} (y_{ijk} - \bar{y}_{+++})^2$	$t_a t_b b - 1$	

Testing Interaction Effects

Idea: exact same as testing treatment effect

1. Biased estimator of var. if H_0 : no interaction b/w factors is true

$$MS_{\text{int}} = \frac{SS_{\text{int}}}{(t_a - 1)(t_b - 1)}$$

2. Unbiased estimator for variance

$$MS_{\text{error}} = \frac{SS_{\text{error}}}{(t_a t_b - 1)(b - 1)}$$

3. F-statistic

$$F = \frac{MS_{\text{int}}}{MS_{\text{error}}} \sim F_{df_{\text{int}}, df_{\text{error}}}$$

Important note: can only test treatment if no interaction effect b/c interaction implies that treatment effect changes based on another factor

Testing treatments:

$$\frac{MS_A}{MS_E} \sim F_{df_A, df_{\text{error}}}$$

↑

$$H_0: \alpha_1 = \dots = \alpha_{t_A} = 0$$

OR

$$\frac{MS_B}{MS_E} \sim F_{df_B, df_{\text{error}}}$$

↑

$$H_0: \beta_1 = \dots = \beta_{t_B} = 0$$

Advantages of factorial structure

Main advantage: enables discussion on interaction effects & more granular levels on factors.

- Having only treatment ANOVA can't let you conclude if factor levels have diff. b/w each other

For post-hoc analysis, we often use regular treatment structure

Overall interaction analysis process

1. If no treatment effect, finish

2. Look for interaction. If there is, interpret & finish

3. Look at factor effects

4. Proc model assumptions are correct

SAMPLE SURVEY ISSUES

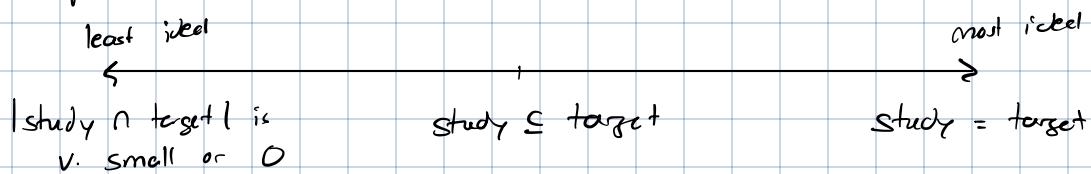
Population is finite → adjust sampling analysis

Terminology

① Sampling unit: individual element in sample which we measure

② Target pop.: pop. we want to study & generalize results

③ Study pop.: pop. that we select samples from

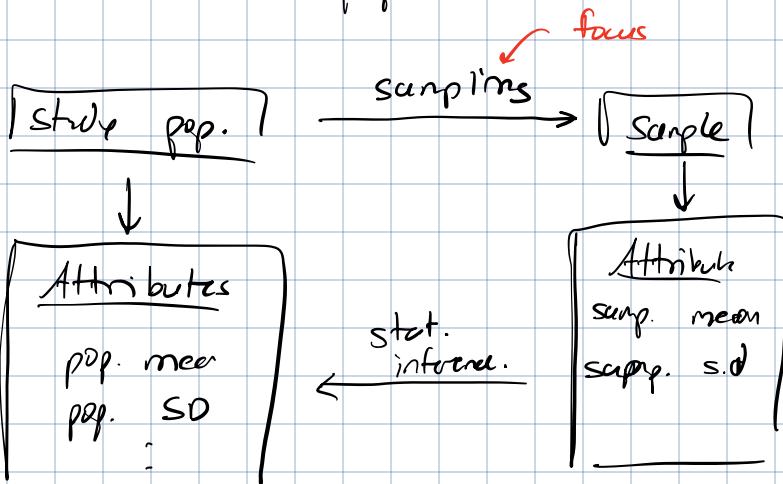


④ Frame: list of units in study pop.

⑤ Sampling: process of gathering data from study pop.

⑥ Attribute: function of pop. that we are interested in (eg. mean, variance, count)

⑦ Census: every unit of pop. is included



Sampling Protocols

Defn: how to select units from study pop.

Sampling protocols

Probabilistic

Assign prob. to A subset
of frame → select sample
based off distn

Benefit: stat. framework to
capture uncertainty
due to sampling

& stat. inference

Non-probabilistic

Subjective selection
of sample

- Eg://
 - Convenience sampling: take what you get
 - Self-selection: units choose to be in sample (eg:// poll)
 - Quota: select units if they match target pop. attr.
 - Judgment: selected s.t. sample is rep. of target pop. off judgment
- No statistical framework to map sample attr. → pop. attr.

Sampling survey errors

- ① Study error: $| \text{target pop. attr.} - \text{study pop. attr.} |$, also called frame error
 - ↳ usually negligible
- ② Sample error: $| \text{study pop. attr.} - \text{sample attr.} |$
 - ↳ Eg:// MAR (missing at random): missing data indep. of value of var. of interest
 - ↳ No bias issues
 - MNAR (missing not at random): \parallel depends on \parallel
 - ↳ Bias issues
- ③ Measurement error: $| \text{true sample unit attr.} - \text{measured sample unit attr.} |$
 - ↳ Caused by mistakes, poor technique, lack of precision

Should be able to determine sources of error given scenario.

PROBABILITY SAMPLING

Major advantage: understand error mathematically

Notation:

N: # of units in frame / pop. (in this course)

U : frame

S : Sample ($S \subset U$)

n : sample size

y_i : response variable of unit i

Sampling protocol

Defn: how to choose samples. Uniquely det. by $P(s)$ / steps to choose s

$$\hookrightarrow \text{Constraint: } \sum_{S \in D} P(S) = 1$$

\hookrightarrow All possible samples

Ex:// $U = \{1, 2, 3\}$. Potential samples $D = \{S_1, S_2, S_3, S_{12}, S_{13}, S_{23}, S_{123}\}$

Protocol A: $P(S_{12}) = P(S_{23}) = P(S_{13}) = \frac{1}{3}$

All other: $P(S) = 0$

$$P(S) = \begin{cases} \frac{1}{3} & \text{if } |S| = 2 \\ 0 & \text{o.w.} \end{cases}$$

Protocol B: $P(S_1) = P(S_{12}) = P(S_{23}) = P(S_{13}) = \frac{1}{4}$

$P(S) (S \neq \emptyset) = 0$.

Inclusion & joint probabilities

Inclusion prob: $P_i = P(i \in S)$

{ Does not define protocol.

Joint incl. prob: $P_{ij} = P(i \in S \wedge j \in S)$

Ex://

s	s_1	s_2	s_3	s_{12}	s_{13}	s_{23}	s_{123}
Sample	{1}	{2}	{3}	{1, 2}	{1, 3}	{2, 3}	{1, 2, 3}
$P(s)$	1/4	0	0	1/4	1/4	0	1/4

$$P_1 = P(s_1) + P(s_{12}) + P(s_{13}) + P(s_{123}) = 1$$

$$P_2 = P(s_2) + P(s_{12}) + P(s_{23}) + P(s_{123}) = \frac{1}{2}$$

$$P_{12} = P(s_{12}) + P(s_{123}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$$

$$P_{23} = P(s_{23}) + P(s_{123}) = \frac{1}{4}$$

Common sampling protocols

① Simple random sampling w/out replacement (SRSWOR)

Design: every sample of size n has equal chance

Intuition: 1 unit from pop. \rightarrow sample until you have n units.

$$P(s) = \begin{cases} \frac{1}{\binom{N}{n}} & |s|=n \\ 0 & \text{o.w.} \end{cases}$$

$$P_i = \frac{\binom{N-1}{n-1}}{\binom{N}{n}} = \frac{n}{N}$$

$$P_{ij} = \frac{\binom{N-2}{n-2}}{\binom{N}{n}} = \frac{n(n-1)}{N(N-1)}$$

② Stratified random sampling

Idea: sep. pop. into strata (variation within strata \downarrow , b/w strata var. \uparrow)

then choose n_h units from strata h , $\sum n_h = n \Rightarrow$ sample

$$P(s) = \begin{cases} \frac{1}{\sum_{i=0}^h \binom{N_i}{n_i}} & |s|=n \\ 0 & \text{o.w.} \end{cases}$$

$$P_i = \frac{\binom{N_{i-1}}{n_{i-1}} \sum_{j=0, j \neq i}^h \binom{N_j}{n_j}}{\sum_{i=0}^h \binom{N_i}{n_i}}$$

③ Cluster sampling

Idea: sep. pop. int. clusters (var. in cluster ↑, var. b/w clusters ↓)

Randomly sample clusters → census on each chosen cluster

④ Systematic sampling

Idea: Choose every m^{th} unit in pop., incmt m , do again → each time is new sample.

⑤ Two-stage sampling

First stage: cluster sample. 2nd stage: SRSWOR / cluster chosen.

Population parameters

Example params. of interest

$$\mu = \frac{1}{N} \sum y_i \quad \text{ans.} \quad \gamma = \sum y_i \quad \text{total.} \quad \sigma^2 = \frac{1}{N-1} \sum (y_i - \mu)^2 \quad \text{variance.}$$

Proportion:

Assume variable of interest is on indicator r : $z_i = \begin{cases} 1 & \text{--} \\ 0 & \text{--} \end{cases}$

$$\mu_z = \frac{\sum z_i}{N} = \frac{\pi}{N} = \bar{\pi} \Rightarrow \text{proportion is cons. of indicator var.}$$

Properties of variance:

$$\textcircled{1} \quad \sigma^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \mu)^2 = \frac{1}{N-1} \left\{ \sum_{i=1}^N y_i^2 - N\mu^2 \right\}$$

$$\textcircled{2} \quad \sigma_z^2 \approx \pi(1-\pi) \quad \text{as } N \rightarrow \infty$$

Proof:

$$\sigma_z^2 = \frac{1}{N-1} \sum_{i=1}^N (z_i - \mu_z)^2$$

$$= \frac{1}{N-1} \left(\sum_{i=1}^N z_i^2 - N\mu_z^2 \right) \quad \text{①} \quad | \quad z_i = \begin{cases} 1 & \text{--} \\ 0 & \text{--} \end{cases}$$

$$= \frac{1}{N-1} \left(\sum_{i=1}^N z_i - N\mu_2^2 \right)$$

$$z_i^2 = \begin{cases} 1 & i \in S \\ 0 & \text{o.w.} \end{cases}$$

$$= \frac{1}{N-1} (N\mu_2 - N\mu_2^2)$$

$$= \frac{N(\mu_2 - \mu_2^2)}{N-1}$$

$$= \frac{N}{N-1} \pi(1-\pi)$$

$$\therefore \lim_{N \rightarrow \infty} \sigma_z^2 = \pi(1-\pi)$$

Estimators of param.

$$\hat{\mu} = \frac{1}{n} \sum y_i, \quad \hat{\sigma}^2 = \frac{1}{n-1} \sum (y_i - \hat{\mu})^2, \quad \hat{Z} = N\hat{\mu}$$

Estimators of params:

$$\hat{\mu} = \frac{1}{n} \sum_{i \in S} y_i \quad \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i \in S} (y_i - \bar{y})^2$$

randomness
↓
def. by protocol.

$$E[\text{R.V.}] = \sum_{s \in D} (\text{R.V. value}) \cdot P(s)$$

Simple random sampling w/out replacement (SRSWOR)

Three results:

$$\textcircled{1} \quad E[\hat{\mu}] = \mu$$

Proof:

$$E[\hat{\mu}] = E \left[\frac{1}{n} \sum_{i \in S} y_i \right]$$

$$= E \left[\frac{1}{n} \sum_{i=1}^n y_i \cdot I_i \right] \Rightarrow I_i = \begin{cases} 1 & i \in S \\ 0 & \text{o.w.} \end{cases}$$

$$\begin{aligned}
 &= \frac{1}{n} \sum_{i=1}^N y_i \in \{I_i\} \\
 &= \frac{1}{n} \sum_{i=1}^N y_i \cdot P(i \in S) \\
 &= \frac{1}{n} \cdot \frac{n}{N} \sum_{i=1}^N y_i \\
 &= \frac{1}{N} \sum_{i=1}^N y_i \\
 &= \mu
 \end{aligned}$$

$$\textcircled{2} \quad \text{Var}(\tilde{\mu}) = \left(1 - \frac{n}{N}\right) \frac{\sigma^2}{n}$$

finite pop. correction $\rightarrow N \rightarrow \infty, \text{Var}(\tilde{\mu}) \rightarrow \frac{\sigma^2}{n}$

Proof:

$$\begin{aligned}
 \text{Cov}(I_i, I_j) &= E[I_i I_j] - E[I_i] E[I_j] \\
 &= P_{ij} - P_i P_j \\
 &= \frac{n(n-1)}{N(N-1)} \cdot \left(\frac{n}{N}\right)^2 \\
 &= -\frac{n}{N} \cdot \frac{1}{N-1} \left(1 - \frac{n}{N}\right)
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(\tilde{\mu}) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n y_i | I_i\right) \\
 &= \frac{1}{n^2} \left[\sum_{i=1}^n y_i^2 \text{Var}(I_i) + \sum_{i=1}^n \sum_{j \neq i}^n y_i y_j \text{Cov}(I_i, I_j) \right] \\
 &= \frac{1}{n^2} \left[\sum_{i=1}^n y_i^2 \cdot \frac{n}{N} \left(1 - \frac{n}{N}\right) - \sum_{i=1}^n \sum_{j=1}^n y_i y_j \frac{n}{N} \cdot \frac{1}{N-1} \left(1 - \frac{n}{N}\right) \right] \\
 &= \frac{1}{n^2} \cdot \frac{n}{N} \cdot \left(1 - \frac{n}{N}\right) \left[\sum_{i=1}^n y_i^2 - \sum_{i=1}^n \sum_{j \neq i}^n \frac{y_i y_j}{N-1} \right] \\
 &= \frac{1}{n^2} \cdot \frac{n}{N} \cdot \left(1 - \frac{n}{N}\right) \cdot N \sigma^2 \\
 &= \left(1 - \frac{n}{N}\right) \frac{\sigma^2}{n}
 \end{aligned}$$

$$\textcircled{1} \quad E[\tilde{\sigma}^2] = \sigma^2$$

Do proof

Can derive:

$$\begin{aligned} \text{Var}(\hat{\mu}) &= \left(1 - \frac{n}{N}\right) \frac{\sigma^2}{n} \\ \downarrow \\ \text{Standard error: } \sqrt{\text{Var}(\hat{\mu})} \end{aligned}$$

} Useful for inference about μ

C.I. for μ :

$$\frac{\hat{\mu} - \mu}{\sqrt{(1-f) \frac{\sigma^2}{n}}} \sim N(0, 1) \text{ if } N, n \text{ & } (N-n) \text{ are large by CLT}$$

$$\Rightarrow \hat{\mu} \pm c_{\alpha} \text{se}(\hat{\mu}) = \hat{\mu} \pm c_{\alpha} \sqrt{(1-f) \frac{\sigma^2}{n}} \text{ is } 100(1-\alpha)\% \text{ C.I.}$$

Estimation of total:

$$\tilde{\tau} = N \cdot \hat{\mu}$$

Thus:

$$E[\tilde{\tau}] = N \times E[\hat{\mu}] = N\mu$$

$$\text{Var}(\tilde{\tau})^2 = N^2 \times \text{Var}(\hat{\mu}) \Rightarrow \text{SE}(\tilde{\tau}) = N \times \text{SE}(\hat{\mu})$$

$$\begin{aligned} \text{C.I. for total} &= N\hat{\mu} \pm c \cdot N \sqrt{(1-f) \frac{\sigma^2}{n}} \\ &= N(C.I. \mu) \Rightarrow \text{multiply endpoints of C.I. } \mu \text{ for C.I. } \tau \end{aligned}$$

Confidence Intervals and Coverage Probability

Interpretation of C.I.: if we repeat the sampling process & calculate 95% C.I.s, we expect 95% of C.I. to contain true param.

↳ NOT: $P(\text{param in C.I.}) = 0.95$. Why? Param is fixed! Cannot add prob. statements on fixed #

For CLT, we want $N, n, (N-n) \sim 30$.

↳ If all of these are not large, we can still proceed but look at coverage

Converse: proportion of C.I. that contain actual param. Should be close to $100(1-\alpha)$

Proportion estimation

If $Y = \{0, 1\}$, then $\text{Var}(Y) = \pi(1-\pi)$, π is mean. (if $\frac{N}{n} \approx 1$)

$$\therefore \hat{\text{Var}}(Y) = \hat{\sigma}^2 = \hat{\pi}(1-\hat{\pi})$$

$$\therefore \text{S.E.} = \sqrt{\left(1 - \frac{n}{N}\right)} \cdot \frac{\hat{\sigma}}{\sqrt{n}} = \sqrt{\left(1 - \frac{n}{N}\right)} \cdot \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}}$$

By CLT (since $E(Y) = \pi$), we have:

$$\text{C.I. proportion} = \hat{\pi} \pm c \sqrt{\left(1 - \frac{n}{N}\right)} \cdot \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}}$$

If calculating totals of proportions, same trick of multiplying C.I. endpoints exist.

Sample size determination

We often want a C.I. of a certain length \rightarrow need to determine what is optimal sample size.

Just have to use math: $2 \cdot c \cdot \text{SE}(\text{param}) \leq \text{desired size}$.

Ex:// We have $n=50$ & $N=362$. How many units should be sampled if 98% C.I. mean should be no larger than 10?

$$\begin{aligned} \therefore 2 \cdot c \cdot \text{SE}(\bar{\mu}) &\leq 10 \\ 2 \cdot 2.33 \cdot \sqrt{1 - \frac{n}{362}} \cdot \frac{\hat{\sigma}}{\sqrt{n}} &\leq 10 \\ 2 \cdot 33 \sqrt{\frac{1}{n} \cdot \frac{1}{362}} \cdot 17.148 &\leq 5 \\ n &\geq 54.26 \end{aligned}$$

Sometimes, we run a pilot study to get data & estimate params
 ↓
 Sample more to get n.

Should sample at least 55 units.

Ex:// In prev. example, how many units should be sampled if total should be within 500 of true value 95% of the time.

$$c \cdot \text{SE}(\tau) \leq 500$$

$$1.96 \cdot \hat{\sigma} \cdot \sqrt{\frac{1}{n} - \frac{1}{362}} \leq 500$$

$$n \geq 224.62$$

We need sample total of 225, i.e. additional $225 - 80 = 175$ for req. precision.

pilot

Ex:// How many units should be sampled if we want to be 95% confident that proportion is within 0.02 of true proportion?

$$C. S.E.(\hat{\pi}) \leq 0.02$$

$$1.96 \sqrt{\frac{1}{n} - \frac{1}{362}} \cdot \sqrt{0.28(1-0.28)} \leq 0.02$$

$$n \geq 304.97 \Rightarrow \text{Size should be at least 305.}$$

Ex:// N= 1000 chips. Want to estimate proportion of defective chips within 5% of the proportion 19 times out of 20. How large should sample be for this precision

$$C. S.E.(\hat{\pi}) \leq 0.05$$

$$1.96 \sqrt{1 - \frac{n}{N}} \cdot \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}} \leq 0.05$$

$$n \geq \frac{1}{\frac{1}{1000} + \left\lceil \frac{0.05}{1.96 \sqrt{0.25}} \right\rceil^2}$$

We don't know what $\hat{\pi}$ is, so choose most conservative option.
We want n to be as big as possible, so denom should be as small as possible, so $\hat{\pi}(1-\hat{\pi})$ should be as big as possible.

This is when $\hat{\pi}(1-\hat{\pi}) = 0.25$ ($\hat{\pi} = 0.5$)

$$\therefore n \geq \frac{1}{\frac{1}{1000} + \left\lceil \frac{0.05}{1.96 \sqrt{0.25}} \right\rceil^2}$$

$$\geq 277.5$$

278 units should be selected.

RATIO AND REGRESSION ESTIMATING IN SRS

Estimating a Ratio

Goal: Want to estimate ratios (proportion of avg./totals of 2 sub-pops.)

Ex:// Want to estimate grocery expenses / person across households, took a sample of 10 households in Waterloo

Household # (i)	1	2	3	4	5	6	7	8	9	10
Y Last week's Groceries (\$)	200	120	400	80	90	140	210	130	220	60
X # of individuals in the household	3	3	6	2	2	3	4	3	5	1

a) Frame?

List of addresses in Waterloo

b) Parameters:

$$\theta = \frac{\sum_{i=1}^N y_i}{\sum_{i=1}^N x_i} = \frac{\frac{1}{N} \sum_{i=1}^N y_i}{\frac{1}{N} \sum_{i=1}^N x_i} = \frac{\mu_y}{\mu_x} \Rightarrow \text{Both of these are unknown & change blur samples}$$

c) Estimate:

$$\hat{\theta} = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i} = \frac{\bar{y}}{\bar{x}} = \frac{200 + \dots + 60}{3 + \dots + 1}$$

Note that:

Param:	Estimate	Estimator:
$\theta = \frac{\mu_y}{\mu_x}$	$\hat{\theta} = \frac{\bar{y}}{\bar{x}}$	$\tilde{\theta} = \frac{\bar{Y}}{\bar{X}}$

To find distn. of $\tilde{\theta}$, we will linearize via Taylor series

- Recall that linear approx. of $f(x, y)$ about (x_0, y_0) is:

$$f(x, y) \approx f(x_0, y_0) + \left. \frac{\partial f(x, y)}{\partial x} \right|_{(x_0, y_0)} (x - x_0) + \left. \frac{\partial f(x, y)}{\partial y} \right|_{(x_0, y_0)} (y - y_0)$$

- Denning linear approx. for $\tilde{\theta}$:

$$\begin{aligned} \frac{y}{x} &\approx \frac{y_0}{x_0} + \left(-\frac{y_0}{x_0^2} \right) (x - x_0) + \left(\frac{1}{x_0} \right) (y - y_0) \\ \Rightarrow \tilde{\theta} &= \frac{\bar{Y}}{\bar{X}} = \frac{\mu_y}{\mu_x} + \frac{1}{\mu_x} (\bar{Y} - \mu_y) - \frac{\mu_y}{\mu_x^2} (\bar{X} - \mu_x) \end{aligned} \quad \begin{array}{l} y \rightarrow \bar{Y} \\ x \rightarrow \bar{X} \end{array}$$

- Mean:

$$E(\tilde{\theta}) \approx \frac{\mu_y}{\mu_x} + \frac{1}{\mu_x} \left(\cancel{\mu_y} - \mu_y \right) - \frac{\mu_y}{\mu_x^2} \left(\cancel{\mu_x} - \mu_x \right) = \frac{\mu_y}{\mu_x}$$

$\therefore \tilde{\theta}$ is approx. unbiased

• Variance:

$$\begin{aligned}
 \text{Var}(\hat{\theta}) &\approx \text{Var}\left(\frac{\mu_y}{\mu_x} + \frac{1}{\mu_x}(\hat{\mu}_y - \mu_y) - \frac{\mu_y}{\mu_x^2}(\hat{\mu}_x - \mu_x)\right) \\
 &= \text{Var}\left(\frac{1}{\mu_x}\hat{\mu}_y - \frac{\mu_y}{\mu_x^2}\hat{\mu}_x\right) \\
 &= \frac{1}{\mu_x^2} \text{Var}(\hat{\mu}_y - \frac{\mu_y}{\mu_x}\hat{\mu}_x) \\
 &= \frac{1}{\mu_x^2} \text{Var}(\hat{\mu}_y - \theta\hat{\mu}_x)
 \end{aligned}$$

Note that: $\hat{\mu}_y - \theta\hat{\mu}_x = \frac{\sum_{i \in s} y_i}{n} - \theta \frac{\sum_{i \in s} x_i}{n}$

$$\begin{aligned}
 &= \frac{1}{n} \sum_{i \in s} (\underbrace{y_i - \theta x_i}_{r_i}) \\
 &= \frac{1}{n} \sum_{i \in s} r_i \\
 &= \bar{r}
 \end{aligned}$$

$\theta x_i = \frac{\mu_y}{\mu_x} \cdot x_i \approx \hat{y}_i$

B/c this is SRSWOR:

$$\text{Var}(\bar{r}) = \left(1 - \frac{n}{N}\right) \frac{\sigma_r^2}{n} \Leftrightarrow \sigma_r^2 = \frac{\sum_{i=1}^n (r_i - \mu_r)^2}{N-1}$$

$$\therefore \text{Var}(\hat{\theta}) \approx \frac{1}{\mu_x^2} \left(1 - \frac{n}{N}\right) \frac{\sigma_r^2}{n}$$

Unknown

• Estimate of σ_r^2

$$\begin{aligned}
 \hat{\sigma}_r^2 &= \frac{1}{n-1} \sum_{i=1}^n (y_i - \hat{\theta}x_i)^2 \\
 &= \frac{1}{n-1} \sum_{i=1}^n (y_i^2 + \hat{\theta}^2 x_i^2 - 2\hat{\theta}x_i y_i) \\
 &= \frac{1}{n-1} \left(\sum_{i \in s} y_i^2 + \hat{\theta}^2 \sum_{i \in s} x_i^2 - 2\hat{\theta} \sum_{i \in s} y_i x_i \right)
 \end{aligned}$$

Since $\hat{\theta}$ is approx. linear in $\hat{\mu}_y$ & $\hat{\mu}_x \Rightarrow \hat{\theta} \sim N$

Confidence interval:

$$\frac{\hat{\mu}_y}{\hat{\mu}_x} \pm c \cdot \frac{1}{|\mu_x|} \sqrt{1 - \frac{n}{N}} - \frac{\hat{\sigma}_r}{\sqrt{n}}$$

Ex:// Create 95% C.I. for prv. interval

① Get point estimates:

$$\hat{\theta} = \frac{\bar{y}}{\bar{x}} = \frac{\sum y_i}{\sum x_i} = 51.60 \quad \hat{\mu}_x = \bar{x} = 3.2$$

$$\hat{\sigma}_r^2 = \frac{1}{n-1} \left(\sum y_i^2 + \hat{\theta}^2 \sum x_i^2 - 2\hat{\theta} \sum x_i y_i \right) = \frac{1457.6992}{10-1}$$

② Combi:

$$\begin{aligned} \text{C.I. : } \hat{\theta} &\pm 1.96 \cdot \frac{1}{|\hat{\mu}_x|} \sqrt{1 - \frac{n}{N}} \cdot \frac{\hat{\sigma}_r}{\sqrt{n}} \\ &\Rightarrow 51.56 \pm 1.96 \cdot \frac{1}{3.2} \sqrt{1 - \frac{10}{3000}} \cdot \frac{1}{\sqrt{10}} \cdot \sqrt{\frac{1457.6992}{10-1}} \\ &\Rightarrow [43.78, 59.34] \end{aligned}$$

Note that this approach works for looking at ratios of π as well, simply replace.

Ratio estimation of the mean

Idea: We want to estimate μ_y , but if X known & X & Y are corr., can make the estimation of μ_y better?

↳ Only works if $\text{Corr}(X, Y) > 0$

↳ 3 cond. must be met:

① Can measure x in sample

② μ_x (pop.) is known

③ $\text{Corr}(X, Y) > 0$

Mathematically, if we know diff b/w μ_x & $\hat{\mu}_x$, we can correct $\hat{\mu}_y$.

Derivation:

$$\text{Assume } y_i = \theta x_i + \epsilon_i \Rightarrow \mu_y = \theta \mu_x \Rightarrow \theta = \frac{\mu_y}{\mu_x}$$

$$\text{Define: } \hat{\mu}_{\text{ratio}} = \hat{\theta} \mu_x = \frac{\hat{\mu}_y}{\hat{\mu}_x} \mu_x = \hat{\mu}_y \cdot \left[\frac{\mu_x}{\hat{\mu}_x} \right]$$

Corrected!

If X & Y are neg. corr., then $\hat{\mu}_y$ modification would make no sense.

Ex:// Want to estimate brain size (y) w.r.t. weight (x).

$$\mu_x = 172.5, \hat{\mu}_x = 151.1, \hat{\mu}_y = 906784.2$$

$$\therefore \hat{\mu}_{\text{ratio}}(y) = \frac{\bar{y}}{\bar{x}} \cdot \mu_x = \frac{\mu_y}{\mu_x} \cdot \frac{\mu_x}{\bar{x}} = \frac{172.5}{151.1} \cdot 906784.2 = 1085166$$

Sample weight was low \rightarrow inflation needed!

Mean of ratio:

$$\begin{aligned} E[\tilde{\mu}_{\text{ratio}}(y)] &= E\left[\frac{\tilde{\mu}_y}{\tilde{\mu}_x} \cdot \mu_x\right] \\ &= \mu_x \cdot E[\tilde{\theta}] \\ &\approx \mu_x \cdot \theta \quad \text{Taylor approx.} \\ &= \mu_x \cdot \frac{\mu_y}{\mu_x} \\ &= \mu_y \Rightarrow \text{unbiased} \end{aligned}$$

Variance of ratio

$$\begin{aligned} \text{Var}(\tilde{\mu}_{\text{ratio}}(y)) &= \text{Var}\left(\frac{\tilde{\mu}_y}{\tilde{\mu}_x} \cdot \mu_x\right) \\ &= \mu_x^2 \cdot \text{Var}(\tilde{\theta}) \\ &\approx \mu_x^2 \cdot \frac{1}{\mu_x^2} \cdot \text{Var}(\hat{\mu}_y - \theta \hat{\mu}_x) \\ &= \text{Var}(\hat{\mu}_y - \theta \hat{\mu}_x) \leftarrow \sigma_r^2 \\ &= \left(1 - \frac{n}{N}\right) \cdot \frac{\sigma_r^2}{n}, \quad \sigma_r^2 = \frac{\sum_{i=1}^n (r_i - \mu_r)^2}{n-1} \end{aligned}$$

C.I.:

$$\text{C.I.} = \hat{\mu}_{\text{ratio}} \pm c \cdot \sqrt{\left(1 - \frac{n}{N}\right) \cdot \frac{\sigma_r^2}{n}}$$

When is the ratio estimate more precise than the sample avg.?

$$\text{Var}(\hat{\mu}_{\text{ratio}}) < \text{Var}(\bar{y})$$

$$\Leftrightarrow \sum_{i \in S} (y_i - \theta x_i)^2 < \sum_{i \in S} (y_i - \bar{y})^2 \Rightarrow \text{use } \hat{\mu}_{\text{ratio}}$$

- Another key is if the line though the origin of y vs. x explains some of the variation. Why? Model of ratio is $y = \theta x_i + \varepsilon_i$, no intercept!

We could also use LS to estimate θ if $y_i = \theta x_i + \varepsilon_i$

$$\therefore \hat{\theta}_{LS} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$$

Regression Estimation of Mean

Ratio estimate better than sample avg. if line of Y vs. X goes through origin
 ↓
 less variance

If line doesn't go through origin \rightarrow regression estimate needed

Regression model in use:

$$Y_i = \alpha + \beta(x_i - \bar{x}) + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2)$$

From least squares, we have:

$$\hat{\alpha} = \bar{y}, \quad \hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2}$$

Develop regression estimate:

$$\hat{\mu}_{reg} = \hat{\mu}_y + \hat{\beta}(\bar{x} - \hat{\mu}_x)$$

pos. obs. ↓ sample avg. ↙
 corrected $\hat{\mu}_y$ ↑ corrected

- $\hat{\beta} > 0 : \text{larger values of } x \rightarrow \text{larger values of } y$

$\hookrightarrow \bar{x} - \hat{\mu}_x > 0 \Rightarrow \hat{\mu}_{reg} \text{ adj. up}$

$\hookrightarrow \bar{x} - \hat{\mu}_x < 0 \Rightarrow \hat{\mu}_{reg} \text{ adj. down}$

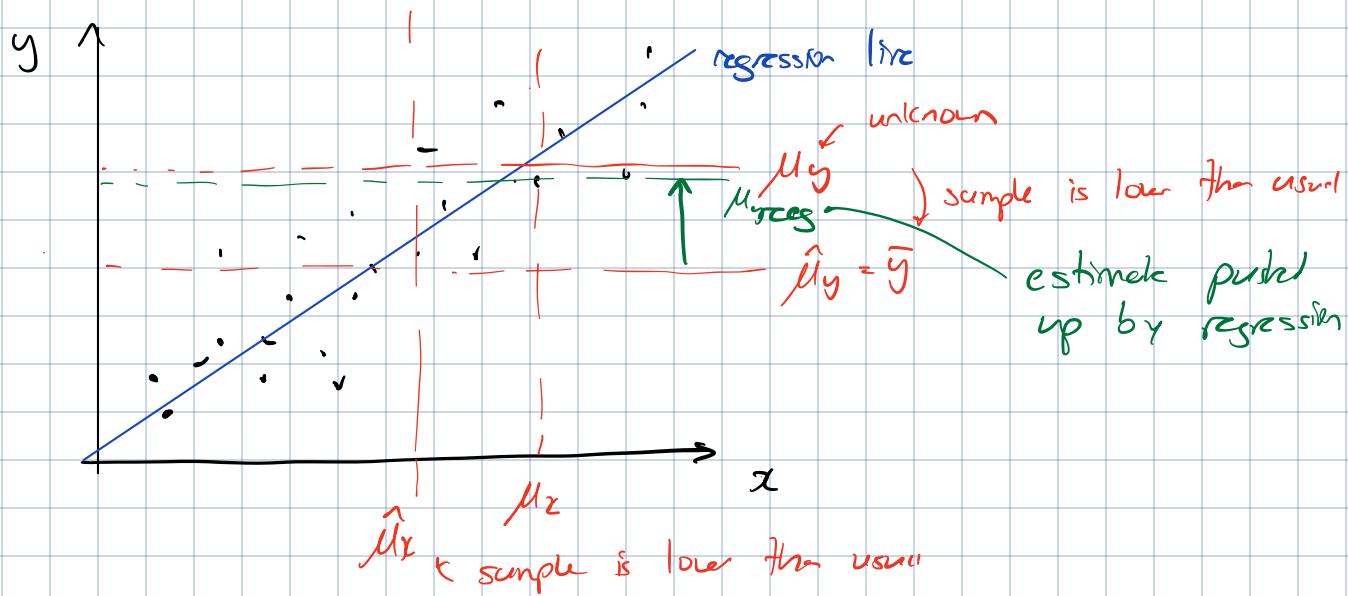
- $\hat{\beta} < 0 : \text{larger values of } x \rightarrow \text{smaller values of } y$

$\hookrightarrow \bar{x} - \hat{\mu}_x > 0 \Rightarrow \hat{\mu}_{reg} \text{ adj. down}$

$\hookrightarrow \bar{x} - \hat{\mu}_x < 0 \Rightarrow \hat{\mu}_{reg} \text{ adj. up}$

Derive from 1st principles of signs.

Ex://



Regression estimator properties

$$\textcircled{1} \quad E[\hat{\mu}_{reg}] \approx \mu_y$$

Proof:

$$\begin{aligned}\hat{\mu}_{reg} &= \bar{y} + \hat{\beta}(\mu_x - \hat{\mu}_x) \\ &= \bar{y} + (\hat{\beta} - \beta + \beta)(\mu_x - \hat{\mu}_x) \\ &= \bar{y} + \beta(\mu_x - \hat{\mu}_x) + (\hat{\beta} - \beta)(\mu_x - \hat{\mu}_x)\end{aligned}$$

$$\begin{aligned}E[\hat{\mu}_{reg}] &= E(\bar{y}) + \beta E(\mu_x - \hat{\mu}_x) + E((\hat{\beta} - \beta)(\mu_x - \hat{\mu}_x)) \\ &= \mu_y + E[(\hat{\beta} - \beta)(\mu_x - \hat{\mu}_x)] \\ &\quad \underbrace{\qquad\qquad\qquad}_{0 \text{ as } n \rightarrow \infty} \quad \underbrace{\qquad\qquad\qquad}_{0 \text{ as } n \rightarrow \infty} \\ &\approx \mu_y\end{aligned}$$

$$\textcircled{2} \quad \text{Var}(\hat{\mu}_{reg}) = \left(1 - \frac{n}{N}\right) \cdot \frac{1}{n} \cdot \boxed{\frac{\sum_{i \in s} (y_i - (\bar{y} + \hat{\beta}(x_i - \bar{x}))^2}{n-1}} \quad \text{SSE}$$

Confidence interval:

$$\hat{\mu}_{reg} \pm c \times \sqrt{\left(1 - \frac{n}{N}\right) \cdot \frac{1}{n} \cdot \frac{\text{SSE}}{n-1}}$$

$\overbrace{\qquad\qquad\qquad}^{\text{S.E. } (\hat{\mu}_{reg})}$

Side note: could we have fit an uncentered model $y_i = \lambda + \beta x_i + \varepsilon_i$?

↳ Yup! β is same, $\hat{\gamma} = \bar{y} - \hat{\beta}\bar{x}$

Centralized model in R: $\text{lm}(Y \sim I(x - \text{mean}(x)))$

Uncentralized : $\text{lm}(Y \sim X)$

Comparison

① Sample avg.

② Ratio estimate

③ Regression estimate

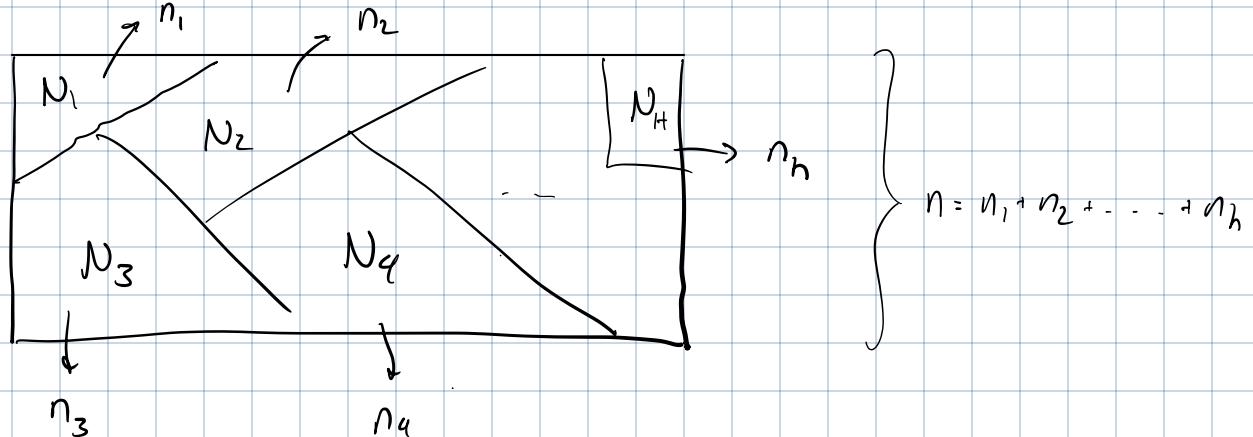
Variance ↑, bias & technicality ↑

STRATIFIED RANDOM SAMPLING

Stratified sampling can yield lower variance if strata chosen carefully

↳ Even though strata is categorical, continuous variables can be bucketed to create strata.

Notation:



μ_h : mean for stratum h , σ_h^2 : variance for stratum h

$$= \frac{1}{N_h-1} \sum_{j=1}^{N_h} (y_{hj} - \mu_h)^2$$

Mean can be considered a weighted avg. of strata:

$$\mu = \frac{N_1}{N} \mu_1 + \dots + \frac{N_h}{N} \mu_h$$

$$= \sum_{h=1}^H w_h \mu_h, \quad w_h = \frac{N_h}{N}$$

① $w_h > 0$

② $\sum w_h = 1$

If we are sampling, then:

$$\textcircled{1} \quad \hat{\mu}_h = \bar{y}_h$$

$$\textcircled{2} \quad \hat{\sigma}_h^2 = \frac{1}{n_{h-1}} \sum (y_{hj} - \hat{\mu}_h)^2$$

$$\textcircled{3} \quad E[\hat{\mu}_h] = \mu_h \quad \& \quad \text{Var}(\hat{\mu}_h) = \left(1 - \frac{n_h}{N_h}\right) \frac{\sigma_h^2}{n_h}$$

Blc SRSWOR applied per stratum.

This leads to our estimate of the pop. mean:

$$\begin{aligned} \hat{\mu}_{\text{strata}} &= \sum_{h=1}^{H} w_h \hat{\mu}_h \\ &= \sum_{h=1}^{H} \frac{N_h}{N} \cdot \hat{\mu}_h \end{aligned}$$

NOT $\frac{n_h}{n}$

Properties:

$$\textcircled{1} \quad E \left[\sum \hat{\mu}_{\text{strata}} \right] = \mu$$

Proof:

$$\begin{aligned} E \left[\sum \hat{\mu}_{\text{strata}} \right] &= E \left(\sum_{h=1}^{H} w_h \hat{\mu}_h \right) \\ &= \sum_{h=1}^{H} w_h \cdot E[\hat{\mu}_h] \\ &= \sum_{h=1}^{H} w_h \cdot \mu_h \\ &= \mu \end{aligned}$$

SRSWOR / strata, so $\hat{\mu}_h$ is unbiased

$$\textcircled{2} \quad \text{Var} \left[\sum \hat{\mu}_{\text{strata}} \right] = \sum_{h=1}^{H} w_h^2 \left(1 - \frac{n_h}{N_h}\right) \cdot \frac{\sigma_h^2}{n_h}$$

Proof:

$$\begin{aligned} \text{Var} \left[\sum \hat{\mu}_{\text{strata}} \right] &= \text{Var} \left(\sum_{h=1}^{H} w_h \hat{\mu}_h \right) \\ &= \sum_{h=1}^{H} w_h^2 \text{Var}(\hat{\mu}_h) + \sum_{h=1}^{H} \sum_{h'=1}^{H} \text{Cov}(w_h \hat{\mu}_h, w_{h'} \hat{\mu}_{h'}) \\ &\quad \text{SRSWOR} \quad \text{Blc non-overlapping strata} \\ &= \sum_{h=1}^{H} w_h^2 \cdot \left(1 - \frac{n_h}{N_h}\right) \frac{\sigma_h^2}{n_h} \end{aligned}$$

Confidence interval

$$\hat{\mu}_{\text{strata}} \pm C \cdot \sqrt{\sum_{h=1}^H w_h^2 \cdot \left(1 - \frac{n_h}{N_h}\right) \frac{\hat{\sigma}_h^2}{n_h}}$$

Important note: estimates for each strata using diff. techniques can be combined

↳ Remember general formula:

$$\hat{\mu}_{\text{strata}} = \sum_{h=1}^H w_h \cdot \hat{\mu}_h, \quad \text{Var}[\hat{\mu}_{\text{strata}}] = \sum_{h=1}^H w_h^2 \cdot \text{Var}(\hat{\mu}_h)$$

Could be simple avg., ratio or regression

Estimating proportion via stratified sampling

Since π_h is just an avg. of an indicator variable:

$$\bar{\pi}_{\text{strata}} = \sum_{h=1}^H w_h \cdot \pi_h, \quad \text{Var}(\bar{\pi}) = \sum_{h=1}^H w_h^2 \left(1 - \frac{n_h}{N_h}\right) \cdot \frac{\pi_h(1-\pi_h)}{n_h}$$

∴ C.I. is:

$$\bar{\pi}_{\text{strata}} \pm C \cdot \sqrt{\text{Var}(\bar{\pi})}$$

Stratified vs. SRSWOR

① Variance of mean estimator

$$\text{Var}(\hat{\mu}_{\text{strata}}) \approx \left(\frac{1}{n}\right) \left(\frac{w_1^2}{w_1} \sigma_1^2 + \dots + \frac{w_h^2}{w_h} \sigma_h^2 \right)$$

$$\frac{w_h}{N} \quad \frac{n_h}{n}$$

Possible that $\text{Var}(\hat{\mu}_{\text{strata}}) > \text{Var}(\hat{\mu})$, but usually σ_i^2 is very small if strata chosen carefully.

② Pop. structure & strata formation

$$\sigma^2 \approx \sum_{h=1}^H w_h \sigma_h^2 + \sum_{h=1}^H w_h (\mu_h - \mu)^2$$

w_h group var. b/w group var.

If strata chosen carefully, $A \downarrow, B \uparrow$. $\text{SE}(\hat{\mu}_{\text{strata}}) \downarrow$ so C.I. narrower

Allocation

① Proportional

Make samples s.t. $\frac{n_h}{n} = \frac{N_h}{N}$ $\rightarrow n_h = w_h \times n$

Then:

$$\begin{aligned}\text{Var}(\hat{\mu}_{\text{strata}}) &= \sum_{h=1}^H w_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{\sigma_h^2}{n_h} \\ &= \sum_{h=1}^H \frac{n_h^2}{n^2} \left(1 - \frac{n_h}{N_h}\right) \frac{\sigma_h^2}{n_h} \\ &= \left(1 - \frac{n}{N}\right) \underbrace{\sum_{h=1}^H w_h \cdot \sigma_h^2}_{\text{Weighted avg. of strata variances.}} \\ &\quad \downarrow \text{smaller than } \sigma^2 \text{ if strata are homogenous}\end{aligned}$$

② Optimal

Choose n_1, \dots, n_h s.t. $\text{Var}(\hat{\mu}_{\text{strata}})$ is minimized.

Minimize via:

$$\min_{n_1, \dots, n_h, n} \left[\sum_{h=1}^H w_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{\sigma_h^2}{n_h} + \lambda \left(\sum_{h=1}^H n_h - n \right) \right]$$

Result:

$$\begin{aligned}n_h &= \frac{w_h \sigma_h}{\sum_{h=1}^H w_h \sigma_h} \times n \\ &= \frac{N_h \sigma_h}{\sum_{h=1}^H N_h \sigma_h} \times n\end{aligned}$$

From this, 2 factors should lead us to sample more from a stratum:

- a) Higher stratum variance
- b) Bigger stratum

Sample Size Determination

We want to find n s.t.

$$C \times \text{S.E.}(\hat{\mu}_{\text{strata}}) \leq \ell$$

We know:

$$\text{Var}(\hat{\mu}_{\text{strata}}) = w_1^2 \left(1 - \frac{n_1}{N_1}\right) \frac{\sigma_1^2}{n_1} + \dots + w_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{\sigma_h^2}{n_h}$$

Assuming optimal allocation:

$$n_h = \frac{w_h \bar{\sigma}_h}{\sum w_h \bar{\sigma}_h} \cdot n$$

$$\text{Assuming } n_i \ll N_i \rightarrow \left(1 - \frac{n_h}{N_h}\right) \approx 1$$

Now:

$$\begin{aligned} \text{Var}(\hat{\mu}_{\text{strata}}) &\approx w_1^2 \cdot \frac{\hat{\sigma}_1^2}{n_1} + \dots + w_h^2 \cdot \frac{\hat{\sigma}_h^2}{n_h} \\ &= \frac{w_1^2 \hat{\sigma}_1^2}{w_1 \hat{\sigma}_1 n} \sum w_h \bar{\sigma}_h + \dots + \frac{w_h^2 \hat{\sigma}_h^2}{w_h \hat{\sigma}_h n} \cdot \sum w_h \bar{\sigma}_h \\ &= \left(\frac{\sum w_h \bar{\sigma}_h}{n} \right) \left(w_1 \hat{\sigma}_1 + \dots + w_h \hat{\sigma}_h \right) \end{aligned}$$

Replacing n_i
w opt. alloc.

Going back to sample size def:

$$C \sqrt{\left(\frac{\sum w_h \bar{\sigma}_h}{n} \right) \left(w_1 \hat{\sigma}_1 + \dots + w_h \hat{\sigma}_h \right)} \leq \ell$$

$$\begin{aligned} C^2 \cdot \frac{\left(\sum w_h \hat{\sigma}_h \right)^2}{n} &\leq \ell^2 \\ n &\geq \left(\frac{C \sum w_h \hat{\sigma}_h}{\ell} \right)^2 \end{aligned}$$

Not variance, std. dev!

From pilot study, estimate $\hat{\sigma}_h$ & calculate n . Use optimal alloc. scheme to get $n_h \forall h$

Ex:// We want to create 95% C.I. of avg. winter hydro usage in kWh, within 50 kWh margin of error.

Stratified households by type of heating system:

	Sample size	Mean	Variance
1 Electricity	24	972	202,396
2 Gas	36	963	96,721

$$N = 35,000 . \quad N_1 = 16,450, \quad N_2 = 18,550$$

a) How much should we sample?

$$\begin{aligned} n &\geq \left(\frac{C \sum U_h \bar{\sigma}_h}{L} \right)^2 \\ &\geq \left[\frac{C}{50} \cdot \left(\frac{16450}{35000} \sqrt{202,396} + \frac{18550}{35000} \sqrt{96721} \right) \right]^2 \\ &\geq 217.6 \end{aligned}$$

We need at least 218 households sampled.

b) How can we allocate to the different strata?

$$\begin{aligned} n_1 &= \frac{1}{\sum U_h \bar{\sigma}_h} \cdot \frac{16450}{35000} \cdot \sqrt{202,396} \cdot 218 \\ &= 122.28 \end{aligned}$$

$$n_2 = 218 - 122 = 96$$

We already sampled 24 & 36 units $\rightarrow n_1^* - n_1 = 122 - 24$
 $n_2^* - n_2 = 96 - 36$

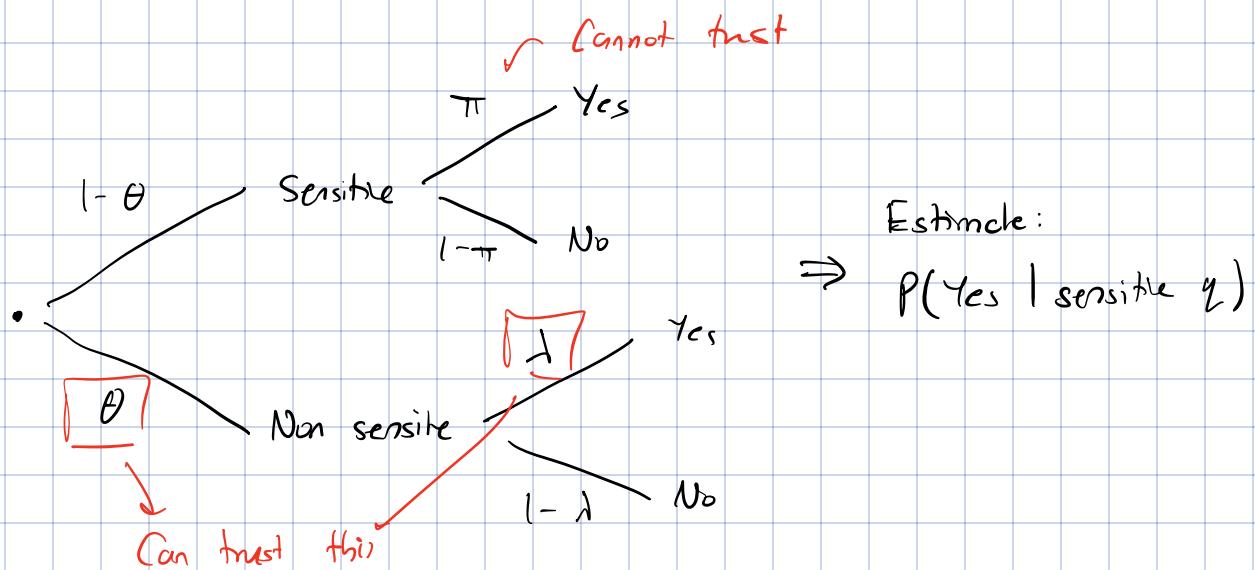
Estimating a Controversial / Sensitive Proportion

Sometimes, surveys ask sensitive question which may yield untruthful responses

To estimate successfully, we have 3 potential methods:

- ① Mix question w/ neutral questions
- ② Let unit answer question @ random (may/may not be sensitive)
- ③ Estimate via total probability rule

Methodology



$$P(\text{Yes}) = P(\text{Yes} \mid \text{Non-sensitive}) + P(\text{Yes} \mid \text{Sensitive})$$

$$\Rightarrow p = \theta\lambda + (1-\theta)\pi$$

$$\therefore \pi = \frac{\hat{p} - \theta\lambda}{1-\theta}$$

$$\Rightarrow \hat{\pi} = \frac{\hat{p} - \theta\lambda}{1-\theta}$$

Deriving C.I. of $\hat{\pi}$:

$$\begin{aligned} \text{Var}(\hat{\pi}) &= \left(\frac{1}{1-\theta}\right)^2 \cdot \text{Var}(\hat{p} - \theta\lambda) \\ &= \left(\frac{1}{1-\theta}\right)^2 \cdot \text{Var}(\hat{p}) \\ &\approx \left(\frac{1}{1-\theta}\right)^2 \cdot \left(1 - \frac{n}{N}\right) \cdot \frac{p(1-p)}{n} \end{aligned}$$

Thus:

$$C.I. \quad \frac{\hat{p} - \theta}{1 - \theta} \pm C_{\alpha} \cdot \sqrt{\left(\frac{1}{1-\theta}\right)^2 \cdot \left(1 - \frac{n}{n}\right) \cdot \frac{\hat{p}(1-\hat{p})}{n}}$$

This method works best if $n \uparrow$, $\theta \downarrow$ (might lead to invalid ans.)

$\hat{\pi}$ may be > 1 or $< 0 \rightarrow$ set to 1 or 0 respectively.