Introduction 560 Hierarchical modeling

Peter Hoff

Statistics, University of Washington

Multilevel data

Multilevel data: Data for which there are

- multiple nested levels of sampling, and/or
- multiple nested sources of variability.

Such data are also often called hierarchical data.

Examples:

Educational testing: students are nested within classes, classes within schools.

Agricultural experiments: subplots nested within whole plots.

Clinical trials: observations over time nested within patients, patients within hospitals.

Terminology

observational unit: an object for which data maybe observed.

macro-level unit: a unit within which other units are nested.

micro-level unit: a unit nested within another unit.

Note that unless there are only two levels, macro versus micro is relative.

Synonyms:

macro units	micro units
top-level units	bottom-level units
primary units	secondary units
clusters	units
groups	units

Notation:

 $y_{i,j}$ = response of *i*th micro unit in *j*th macro unit

Why multilevel data?

Consider the costs of administering an in-person survey to:

- 100 randomly sampled public high-school students in Washington state;
- 10 students sampled from 10 randomly sampled high schools.

Cluster sampling

The second sampling scheme is called *cluster sampling* or *two-stage sampling*.

Cluster sampling

- is often cheaper per sampled unit;
- often gives less reliable estimates of population means.

Task: Estimation of a population mean

Task: Estimate the population mean μ from sample data.

Questions:

- How do cluster sampling and SRS compare?
- How do you infer μ from cluster sample data?



 $\mu = 2.1124814$



 $\mu = 2.1124814$



 $\mu{=}2.1124814$, $\bar{y}{=}1.8687553$



 $\mu = 2.1124814$



 $\mu = 2.1124814$



 $\mu{=}2.1124814$, $\bar{y}{=}2.2758576$

Variability of sample mean





 μ =2.1124814



 $\mu{=}2.1124814$, $\bar{y}{=}2.0593259$



 μ =2.1124814



 $\mu {=} 2.1124814$, $\bar{y} {=} 2.20719$

Variability of sample mean



Comparison of sampling variability



Heterogeneity, homogeneity and dependence

As we will show mathematically,

across-group heterogeneity \Leftrightarrow within-group homogeneity

 $\Leftrightarrow {\sf within-group\ correlation\ or\ dependence}$

Equivalently, across-group heterogeneity generally

- increases the variance of the sample;
- reduces the precision of estimates.

 $Var[\bar{y}_{tss}] \geq Var[\bar{y}_{srs}]$

Task: Construct a 95% CI for the population mean.

t-interval for SRS:

If y_1, \ldots, y_n is an iid sample with $E[y_i] = \mu$ and $Var[y_i] = \sigma^2$,

$$\mathsf{E}[\bar{y}] = \mu \,\,,\,\, \mathsf{Var}[\bar{y}] = \sigma^2/n.$$

By the central limit theorem,

$$ar{y} \stackrel{.}{\sim} \mathsf{N}(\mu, \sigma^2/n) \;,\; rac{ar{y} - \mu}{\sigma/\sqrt{n}} \stackrel{.}{\sim} \mathsf{N}(0, 1).$$

As σ^2 is generally unknown, we use

$$rac{ar y-\mu}{s/\sqrt{n}} \sim t_{n-1}, \ , ext{where} s^2 = rac{1}{n-1}\sum(y_i-ar y)^2.$$

From this, we have

$$ar{y} \pm t_{n-1,.975} imes s/\sqrt{n}$$
 is a 95% CI for μ .



$$\bar{y} \pm t_{n-1,.975} \times s/\sqrt{n}$$

What if we apply the formula to data from a cluster sample? If y_1, \ldots, y_n are from a SRS, then

$$\operatorname{Var}[\bar{y}] = \sigma^2/n = \operatorname{E}[s^2/n].$$

 s/\sqrt{n} provides a good estimate of the sd of \bar{y} .

If y_1, \ldots, y_n are from a cluster sample, then generally

$$\operatorname{Var}[\bar{y}] > \sigma^2/n \approx \operatorname{E}[s^2/n].$$

 s/\sqrt{n} is generally an underestimate of the sd of \bar{y} . How will the resulting confidence interval behave if sd $(\bar{y}) > s/\sqrt{n}$?



Summary:

- Across group heterogeneity = within group similarity.
- Within group similarity leads to positively correlated cluster sample data.
- The variance of the sample mean from (positively) correlated data is higher than that of the mean of uncorrelated data.
- Statistical inference ignoring such correlation will be inaccurate.

Moving forward: We will develop techniques to

- evaluate within and across-group heterogeneity;
- provide accurate statistical inference based on cluster samples.

Task: Estimation of an effect

Suppose

- $x \in \{0, 1\}$
- $\mu_1 = \mathsf{E}[y|x = 1]$
- μ₀ = E[y|x = 0]

Task: Estimate the difference $\delta = \mu_1 - \mu_0$ based on cluster sample data.

Data: For each group j, we have $(y_{1,j}, x_{1,j}), \ldots, (y_{n,j}, x_{n,j})$.

Question: What could go wrong by ignoring the multilevel nature of the data?





Ignoring group differences can lead to overconservative inference.



- The population mean difference is zero;
- The sample mean difference based on any two groups is not zero.

Ignoring group differences can lead to underconservative inference.



• $\mu_1 - \mu_0 > 0$ in the overall population;

• $\mu_{1,j} - \mu_{0,j} < 0$ in every group.

Micro/group effects may be different from macro/population effects.

This is sometimes called *Simpson's paradox*.

Summary:

- Across group heterogeneity can lead to *over or under* conservative inference.
- Aggregated macro effects may be different from micro effects.
- Statistical inference ignoring groups can be inaccurate in *unpredictable* ways.

Moving forward: We will develop techniques to

- differentiate between macro and micro level effects;
- appropriately control for within and between-group heterogeneity.

X, x are macro and micro level explanatory variables

Y, y are macro and micro level outcome variables



What are the effects of SES (*x*) on political opinion (*y*)? (a *micro effect*)

X, x are macro and micro level explanatory variables

Y, y are macro and micro level outcome variables



What are the effects of State GDP (X) on political opinion (y) ? (a macro-micro effect)

X, x are macro and micro level explanatory variables Y, y are macro and micro level outcome variables



What are the effects of State GDP (X) on statewide political opinion (Y)? (a *a macro effect*)

X, x are macro and micro level explanatory variables

Y, y are macro and micro level outcome variables



What are the effects of State GDP (X) and SES (x) on political opinion (y)? (a *multilevel effects*)

Example: Income and voting patterns

Exit poll data from 2004 presidential election

- $j \in \{1, \dots, 50\}$ indexes the states,
- $y_{i,j}$ is the voting variable for person *i* in state *j*,
- $x_{i,j}$ is the measure of income for person (i,j).

Macro-level income-voting relationships



Micro-level income-voting relationships



Joint estimation of effects

In general we may be interested in understanding all of the following:

- macro level effects,
- micro level effects,
- heterogeneity of micro effects across groups.