

## Index of Lecture 10a: Introduction to Clustered Data

Page	Title
1	Practical information
2	What is clustering and clustered data?
3	Hierarchical data structure
4	Other data structures with clustering
5	Somatic cell count datasets
6	Reunion Island study
7	Why is clustering important?
8	Simulated data results
9	Variance inflation
10	Example framework
11	Mixed model for hierarchical data
12	Random effects (recap)
13	Mixed model for somatic cell count data
14	Variance components
15	Estimation in linear mixed models
16	Statistical inference in linear mixed models
17	Statistical inference (continued)

## PRACTICAL INFORMATION

### Today's session:

- introduction to clustered data,<sup>1</sup>
  - \* **explanation of idea/concept** and its relation to data structures,
  - \* **impacts of clustering** on data analysis,
- hierarchical linear mixed models,<sup>2</sup>
  - \* introduction to likelihood-based analysis and interpretation, focusing on main principles,
  - \* logistic mixed models → next session,
  - \* more details in: multilevel summer course (July 14-18, see also announcement at <http://cver.upei.ca>),

### News/Schedule:

- lab session for this week's material on Monday (no homework for Friday),
- home assignment 4 and project proposal due on Monday (March 17).

---

<sup>1</sup> Corresponding to Sections 20.1-3 of the VER2/MER textbook (with some skipings).

<sup>2</sup> Sections 21.1-2 and parts of 21.5 from the VER2/MER textbook.

## WHAT IS CLUSTERING AND CLUSTERED DATA?

This **terminology** is used in the VER2/MER texts and broadly in vet-epi, but is **not unique or standard**, and the term cluster is also used in multivariate and spatial analysis with a somewhat different meaning.<sup>3</sup>

**Clustering** of data (values) loosely represents situations when some observations are **more similar** than others, for reasons **not explained by other factors** of interest or accounted for explicitly in the modelling. **Simple examples:**

- body dimensions of humans/animals are “more similar” for individuals of same age and gender (and breed), but these relations are typically known and controlled for,
- performance is (often) “more similar” of individuals within the same unit (e.g. family, school, workplace, farm) than between different units; although individuals’ membership to units is known, such relations can exist beyond what is explained by unit characteristics ( $\Rightarrow$  is unexplained variation).

**Statistically**, unexplained (dis)similarity between values is expressed by their correlation:

- more (less) similar within a cluster  $\Rightarrow$  positively (negatively) correlated,
- independence  $\sim$  no clustering.

**Key message:** clustering implies **lack of independence!**

<sup>3</sup> Observations within a cluster are closer together (around the cluster centre) than between clusters.

## HIERARCHICAL DATA STRUCTURE

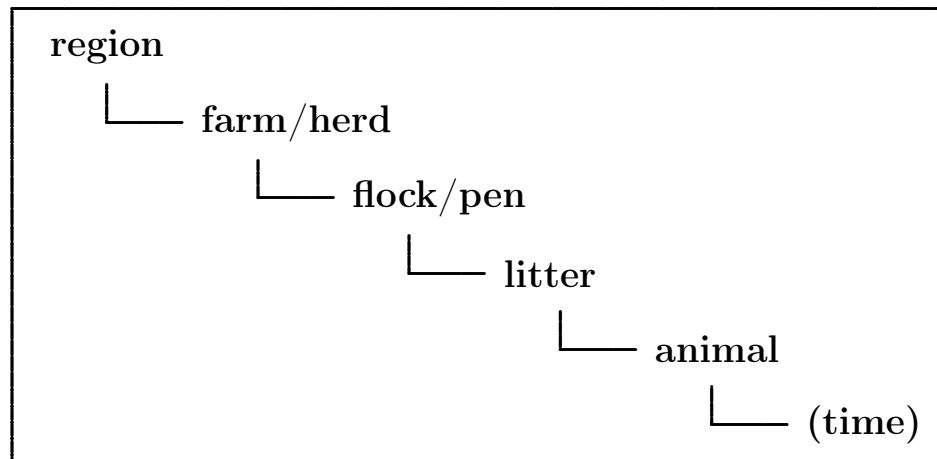
A **hierarchical data structure**, typically implying both of the following:

- observations grouped at different levels,
- factors (e.g. treatments) reside (or are applied) at different levels.

may induce **clustering** in the data, that is, some observations are more alike than others, or put in another way: the **observations are no longer independent**.

**Typical example** from  
veterinary epidemiology:

(usually, only some of the  
levels are present)



**Note:** “time” as a bottom level  $\sim$  longitudinal data / repeated measures on the same animal, and raises some extra modelling issues (next slide, and Session 12 of VHM 802).

A common usage is that levels are **nested** within each other, e.g. that animals are nested within litters or pens are nested within herds.

## OTHER DATA STRUCTURES WITH CLUSTERING

- **Repeated measures:**
  - \* several measurements of a variable taken on the same animal (or other unit of observation) over time,
  - \* on previous slide indicated as an extra level at the bottom of the hierarchical structure diagram, but repeated measures also possible at intermediate levels,
  - \* correlation between two observations may depend on the time between them,
- **Spatial data structure:**
  - \* correlation between two units of observation may depend on the physical distance between them (e.g. herd locations, cows within a barn, or plots in an agricultural field),
  - \* spatial correlation may occur at any level in the hierarchy,
- **Non-nested data structures:**
  - \* may be due to imperfections in a hierarchical structure, e.g. an animal switches herds during the study,
  - \* may also be because levels are **not nested** but **crossed**, whereby units at one level are not necessarily together at higher levels (e.g., fish at the same farm are not necessarily from the same hatchery),
  - \* also **multiple membership**, where a unit “belongs to” multiple higher level units.

## SOMATIC CELL COUNT DATASETS

scc\_40 — a real somatic cell count dataset:

A subset comprising 40 herds from a large dataset collected in 1993-94 by Jens Agger and co-workers including about 2150 Danish herds and 150 000 cows followed throughout one lactation. The data contain approximately monthly milk records plus information collected through herd questionnaires.

Variable	Description	Values
herdid	herd id	1 – 40
cowid	cow id	1 – 2178
test	approximate month of lactation	0 – 10
t_lnscc	natural log scc (in 1000s) on test day	2.3 – 9.2
t_dim	days in milk on test day	10 – 305
t_season	season of test day	1 – 4 (1 = Jan-Mar, etc.)
c_heifer	parity of cow	0/1 (1 = heifer)
h_size	average herd size	10.3 – 101.5
t_ecm	energy-corrected <sup>1</sup> milk yield	2.2 – 68.5

<sup>1</sup> computed by the formula:

$$\text{ecm} = \text{kgmilk}(0.383 \text{ fatpct} + 0.242 \text{ proteinpct} + 0.7832)/3.14.$$

**Subdataset** scc40\_2level: only first observation per lactation included  
 ⇒ one observation per cow.

## REUNION ISLAND STUDY

- carried out 1993-1996 on Reunion Island (Emmanuel Tillard, CIRAD),
- data analyzed 1996-2000 (with a strong helping hand from Ian Dohoo), and results communicated to the cattle industry and published 1999-2001.

**Objective:** identify the factors and levels at which most of the variability in reproductive performance resides ... where interventions are likely to have the most effect.

**Reproductive performance in cows** measured by time from calving to conception, which is composed of

- time from calving to first service,
- conception from first service (yes/no),
- if no, time from first service to conception.

**Data size and structure:**

Level	Number of units	Per unit at level above	
		Average	Range
Region	5	—	—
Herd	50	10.0	3–16
Cow	1575	31.5	8–105
Lactation	3027	1.9	1–5

- no cow movements between herds (~ strict hierarchical structure).

## WHY IS CLUSTERING IMPORTANT?

Three main answers:

- it may have **strong and not easily predictable effects** on the statistical analysis, because it invalidates the assumption of independence made by models such as multiple linear and logistic regression — a few remarks on **ignoring clustering**:
  - \* may affect both the **estimates** themselves (bias!) and their **precision** (confidence intervals, tests),
  - \* the impact depends on the **statistical model/method** used, and on whether the **predictor** varies between or within clusters,
  - \* even when the **clustering is small** and perhaps not significant, it should be accounted for whenever feasible.
- it leads to a **separation of the variation** in the data which may be of intrinsic interest; **examples**:
  - \* **somatic cell count** example (2-level version):  
variation split into between-herd and within-herd (or between-cow) variation,
  - \* **reproductive performance in Reunion Island** example:  
variation split into components for (region), herd, cow and lactation,  
— the components indicate where there is largest potential for improvement,
- cluster units (e.g., farms) may have a role as **confounders**.

## SIMULATED DATA RESULTS

**Simulated datasets** (so we know what the true values are...):

- 100 herds with an average of 116 cows (range: 20–311),
- single binary predictor  $x$ , with equal distribution of 0's and 1's, either at the **herd level** (“feeding type”) or at the **cow level** (“hormone treatment”),
- **continuous outcome** (“milk yield”): mean: 30, effect of  $x$ :  $\beta = 5$ ,
- **binary outcome** (“mastitis”):  $\text{logit}(p)$  mean: -1.4, effect of  $x$ :  $\beta = 0.693$  ( $\sim \text{OR} = 2$ ).

**Results** (of a single simulation of each dataset):

Outcome	Continuous $\sim$ linear models						Binary $\sim$ logistic models			
Approach	Ignore clustering		Mixed model		Herd averages		Ignore clustering		Mixed model	
Parameter	estimate	SE	estimate	SE	estimate	SE	estimate	SE	estimate	SE
	<b>Dataset 1: <math>x</math> as a herd-level factor</b>									
$x$	3.557	0.200	3.796	1.496	3.779	1.497	0.529	0.042	0.620	0.204
constant	30.021	0.146	31.137	1.059	31.166	1.059	-1.242	0.033	-1.305	0.145
	<b>Dataset 2: <math>x</math> as a cow-level factor</b>									
$x$	4.982	0.199	4.968	0.149	–	–	0.586	0.042	0.697	0.046
constant	29.257	0.141	30.646	0.728	–	–	-1.250	0.032	-1.361	0.111

**Conclusion:** effects of clustering seen both in the estimates (especially binary outcome) and their precision (especially herd-level  $x$ ).

## VARIANCE INFLATION

**Special case** for understanding the effect of clustering:

herd (cluster) level predictor  $x$  (quantitative/categorical), herds of equal size ( $m$ ) and (for simplicity) normally distributed outcome:

- **estimate(s) for  $x$** : unaffected by the clustering,
- **SE(s) for  $x$** : multiplied by  $\sqrt{\text{VIF}}$ , the variance inflation factor (VIF<sup>4</sup>), given by<sup>5</sup>

$$\text{VIF} = 1 + (m - 1) \rho,$$

where  $\rho$  is the intra-class correlation (ICC) between two observations in a herd,

- VIF may be interpreted as a measure of the loss of accuracy/power due to the clustering.

**Some implications** (for this particular situation):

- easy to correct for clustering (scale up SEs by  $\sqrt{\text{VIF}}$ ),
- easy to **correct sample size calculations for clustering**: (multiply required sample sizes by VIF),
- a low ICC with a moderate group size can have as much impact as a high ICC with a small group size,
- even low ICCs lead to **high variance inflation** when the group size is large.

<sup>4</sup> Not related to variance inflation factors for collinearity in multiple regression.

<sup>5</sup> The variance inflation is basically on the herd means:  $\text{Var}(\bar{Y}) = (\sigma^2/m) \times \text{VIF}$ .

## EXAMPLE FRAMEWORK

Consider the following problem:

- study of risk factors for (high) somatic cell counts (e.g., as a crude indicator of mastitis),
- **one recording** (for simplicity) of the cell count in a milk sample **from each cow**; in total,  $n$  cows,
- additional recordings of **explanatory variables for each cow**, such as lactation stage (days in milk), age, breed,...
- also **explanatory variables at the herd level**, e.g. housing type, herd size,...
- **linear model**:

$$Y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_k x_{ki} + \varepsilon_i, \quad \text{or} \quad Y = X'\beta + \varepsilon$$

where

- \*  $Y_i$  = (natural) log somatic cell count for cow  $i$ ,  $i = 1, \dots, n$ ,
- \*  $x_{ri}$ 's contain values of the explanatory variables,<sup>6</sup>
- \*  $(\beta_0), \beta_1, \dots, \beta_k$  are regression coefficients for  $x$ 's,
- \*  $\varepsilon_i$  = error term  $\sim N(0, \sigma^2)$ .

---

<sup>6</sup> In this notation (from the VER2 textbook), we use  $x_{ri}$  instead of the usual  $x_{ir}$ ;  $X' = (x_{ri})_{ir}$  is the  $n \times (k+1)$  design matrix, including as  $(x_{0i})$  a column of 1's.

## MIXED MODEL FOR HIERARCHICAL DATA

**Simplest case**  $\sim$  extended cell count example, for measurements on cows in several herds,

$$Y_{ij} = \beta_0 + \beta_1 x_{1ij} + \dots + \beta_k x_{kij} + u_j + \varepsilon_{ij}, \quad \text{or} \quad Y = X'\beta + u + \varepsilon$$

where

- $Y_{ij}$  = log somatic cell count for cow  $i$  in herd  $j$ ,
- $u_j$  = random  $j^{\text{th}}$  herd effect  $\sim N(0, \sigma_h^2)$ ,
- $\sigma_h$  = scale of random herd effects, interpretable as the amount of random variation in log-scc between herds;
  - \* e.g., 95% of herds expected within  $0 \pm 1.96 \sigma_h$ ,
- $i$  = cow number,  $j$  = herd number.

**Definitions** (non-Bayesian terminology):

- “**random**” effect: a model term (right hand side) which is a random variable (often not counting the error term  $\varepsilon$ ),
- “**fixed**” effect: modelled by non-stochastic parameters ( $\beta$ 's, often not counting the intercept  $\beta_0$ ),
- **mixed model**: both fixed and random effects.

## RANDOM EFFECTS (RECAP)

- the only new concept in mixed models,
- enable a **separation (and quantification)** of variation at different levels in the data,
- enable **correct analysis of predictors at different levels** within the same model,
- involve **additional assumption(s)** of normal distribution and variance homogeneity.

### Motivations for random effects (decreasing importance):

- hierarchical data structure:
  - \* **rule**: insert a random effect for each hierarchical level above the bottom level (exceptions: 10aL–17),
- correct analysis of treatments allocated to **larger experimental units** (“split-plot” idea<sup>7</sup>),
- factor where **interest is in the variation between units** (within a level) rather than specific units in study:
  - \* units may be randomly selected,
  - \* units should **represent “population”** – to which the conclusions from the study may be generalized,
- avoid many (nuisance<sup>8</sup>) parameters in model/estimation.

<sup>7</sup> Split-plot designs are experimental designs where treatment factors are applied to units of different sizes; discussed in detail in regular version of VHM 802, and in the GO textbook.

<sup>8</sup> Of no or little intrinsic interest.

## MIXED MODEL FOR SOMATIC CELL COUNT DATA

- o **data**: 2-level somatic cell count data (scc40\_2level),
- o **outcome**: log somatic cell count (t\_lnscc),
- o **fixed effects**: season (categorical), dim, heifer, hsize (all quantitative),
- o **random effects**: herds (note: only 1 observation per cow in these data).

```
. mixed t_lnscc h_size c_heifer i.t_season t_dim || herdid:, reml
```

```
Mixed-effects REML regression          Number of obs      =      2178
Group variable: herdid                 Number of groups   =        40
                                       Obs per group: min =        12
                                       avg =          54.5
                                       max =          105
                                       Wald chi2(6)       =      244.36
Log restricted-likelihood = -3624.9622  Prob > chi2        =      0.0000
```

```
-----+-----
      t_lnscc |      Coef.   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----+-----
      h_size |   .0040837   .0037726     1.08  0.279    - .0033105   .0114778
...
      _cons  |   4.641202   .1974215    23.51  0.000     4.254263   5.028141
-----+-----
```

```
Random-effects Parameters |   Estimate  Std. Err.   [95% Conf. Interval]
-----+-----
herdid: Identity         |
      var(_cons) |   .1491533   .0436191     .0840821   .2645832
-----+-----
      var(Residual) |   1.557228   .0477206     1.466451   1.653625
-----+-----
```

```
LR test vs. linear regression: chibar2(01) =    97.01 Prob >= chibar2 = 0.0000
```

## VARIANCE COMPONENTS

Always a **decomposition of the total variation**:

total variation = fixed effects variation + random variation,<sup>9</sup>

Additionally in random effects models — a **decomposition of the random variation**.

**2-level model** — cell count example (herds—cows):

- $\text{Var}(Y_{ij}) = \sigma^2 + \sigma_h^2$  (= total unexplained random variation),
- variance components  $\sigma^2$  and  $\sigma_h^2$ ,
- of the total random variation,  $\sigma_h^2/(\sigma^2 + \sigma_h^2)$  resides at the herd level, and the rest at the cow level; the former is also termed a **variance partition coefficient** (VPC),
- $\text{ICC}^{10} = \sigma_h^2/(\sigma^2 + \sigma_h^2)$ , often denoted  $\rho$  (as a correlation coefficient).

**Multilevel models** — Reunion example (herds—cows—lactations):

- $\text{Var}(Y_{ijk}) = \sigma^2 + \sigma_c^2 + \sigma_h^2$ ,
- **proportions of variance** (VPCs) at different levels in the obvious way,
- **two ICCs**:<sup>11</sup>

$$\left\{ \begin{array}{l} \text{lactations of same cow : } (\sigma_c^2 + \sigma_h^2)/(\sigma^2 + \sigma_c^2 + \sigma_h^2), \\ \text{lactations in same herd : } \sigma_h^2/(\sigma^2 + \sigma_c^2 + \sigma_h^2) \end{array} \right\},$$

<sup>9</sup> Least squares decomposition:  $\sum(Y_i - \bar{Y})^2 = \sum((X'\hat{\beta})_i - \bar{Y})^2 + \sum(Y_i - (X'\hat{\beta})_i)^2$

<sup>10</sup> Intra-class (or -cluster) correlation coefficient: the correlation between two observations in the same class/cluster. Alternative **interpretations of clustering** as: variation between clusters, or correlation within clusters.

<sup>11</sup> **General formula**: sum of variance components of **common random effects** divided by sum of all variance components.

## ESTIMATION IN LINEAR MIXED MODELS

“Likelihood”-based estimation assuming normal distributions for all random terms:

- REML (restricted maximum likelihood) or (full) ML:
  - \* theoretical properties differ slightly (REML unbiased, ML less variance),
  - \* in practice only minor differences, unless the number of units at a level is small,
  - \* REML estimates agree with ANOVA-type<sup>12</sup> estimates for balanced data,  
— choice between REML and ML is “a matter of taste” (but keep consistent),
- iterative, numerically robust algorithms,
- performs well for both balanced and unbalanced data,
- available in many statistical software packages, however however
  - \* different modelling flexibility,
  - \* different ability to handle large data structures,
- should give the same estimates from different software packages, up to estimation accuracy, despite minor differences in implementations.

---

<sup>12</sup> A classical statistical method for variance component models relies on the ANOVA-table, and constructs estimates and test statistics from the MS-column; performs well only for almost balanced data; discussed in (the standard) VHM 802 course and the GO textbook.

## STATISTICAL INFERENCE IN LINEAR MIXED MODELS

Tests and confidence intervals — approximate and assuming normal distributions for all random terms:

- **Wald statistics** for fixed effects: based on standard errors and estimated correlations between estimates,
  - \* **95% confidence intervals**:  $\hat{\beta}_r \pm 1.96 \times \text{SE}(\hat{\beta}_r)$ ,  
(better to use  $t^* \sim t$ -distribution with suitable df<sup>13</sup>, but similar if df is large),
  - \* simple to compute, for fixed effect parameters usually ok,
- **likelihood-ratio tests**, based on optimal values of likelihood function (more precisely, differences of  $-2 \log L$ ):<sup>14</sup>

$$G^2 = 2(\log L_{\text{full}} - \log L_{\text{red}}) \sim \chi^2(\text{df}),$$

where df = number of parameters being tested equal to zero,

- \* the **only appropriate test** for **variance parameters**,<sup>15</sup>
- **confidence intervals for variance parameters**: not easy, and usually ok to present Stata's approximate intervals from Wald-type procedure (but inference from these can be misleading).

<sup>13</sup> In Stata 14+, options `dfmethod(satterthwaite)` or `dfmethod(kroger)`; same options in Minitab 18+.

<sup>14</sup> Caution for fixed effects parameters and REML estimation: beware to **not use** the restricted likelihood. Stata's `lrtest` command gives a warning note.

<sup>15</sup> Note: the  $P$ -value should be **half** the value from  $\chi^2(\text{df})$  when testing  $H_0 : \sigma_h^2 = 0$ . (Stata does this per default.)

## STATISTICAL INFERENCE (CONTINUED)

### Model-building guidelines:

- **fixed effects**: similar to linear models,<sup>16</sup>
- **random effects**: generally one for each hierarchical level, but some exceptions:
  - \* fixed effects possible, if **number of units is small** and/or **unrepresentative** of a population, and no predictor has variation at that level — this typically occurs at the highest level,
  - \* level may need to be omitted if only very few replications present in the data.

### Model checking:

- now model assumptions (and potential violations) at multiple levels,
- **predictions, residuals and diagnostics at multiple levels**,
- softwares differ in which of these statistics are available; **Stata** gives easy access to
  - \* lowest level residuals (however, not standardized so well) + fitted values  
⇒ usual model checking at lowest level,
  - \* predicted random effects (“BLUPs”) at higher levels ⇒ normality checks can be done easily, and plots against predicted values with some coding,
- **model checking is an informal, exploratory process.**

<sup>16</sup> **Recommendation**: fixed effects model-building should use models with all relevant random effects.