

# Fitting models to stratified samples

Alastair Scott  
(joint work with Chris Wild)

University of Auckland

## Problem

We want to fit a model,  $f(\mathbf{y}|\mathbf{X};\boldsymbol{\theta})$  say, for the distribution of a response,  $\mathbf{y}$ , as a function of explanatory variables,  $\mathbf{X}$ , to data from stratified samples.

- Particularly interested in situations where the probability that a unit is in a particular stratum may depend on the value of  $\mathbf{y}$

## Example 1: Case-control family study

Study of glioma, the most common type of malignant brain tumour (see Wrensch et al (1997)):

- all newly diagnosed cases in a specified time interval in the San Francisco Bay area
- a population-based sample of controls through random digit dialling.
- brain tumour status and covariate information from the original case-control sample and, in addition, from members of their families.

## Example 1

Want to fit a mixed logistic model

$$\text{logit } P(Y_{ij} = 1 \mid b_i) = \beta_0 + \beta_1 ep_{ij} + \beta_2 ca_{ij} + b_i$$

where  $b_i$  is a random family effect with  $b_i \sim N(0, \sigma^2)$ . Here

- $Y$  is a binary indicator of brain cancer status,
- $ep$  is a binary indicator of epilepsy history,
- $ca$  is a binary indicator of other cancers.

## Example 1

- If we only had information from the original case-control probands, we would just fit a standard logistic ignoring the differential sampling rates (with an adjustment to the constant if wanted).
- Natural to try to mimic this and fit a mixed logistic model ignoring the case-control sampling
  - unfortunately, this doesn't work!

## General set-up

- Population of  $N$  units, split into strata  $\mathcal{S}_1, \dots, \mathcal{S}_L$  with  $N_h$  units in  $\mathcal{S}_h$ .
- We draw independent random samples,  $D_h$ , of  $n_h$  units from the  $N_h$  units in  $\mathcal{S}_h$  ( $h = 1, \dots, L$ ) and observe the corresponding  $(\mathbf{y}, \mathbf{X})$  values.
- Data are of the form:  
 $\{(\mathbf{y}_{hj}, \mathbf{X}_{hj}, j \in D_h), N_h; h = 1, \dots, L\}$ .
  - $N_h$ s provide information about stratum probabilities

## General set-up

- No problem if the probability of being in a stratum depends on  $\mathbf{X}$  alone.
- However, if stratum probabilities depend on  $\mathbf{y}$ , then the sampling is informative and cannot be ignored in the analysis.
  - Sometimes (in many retrospective medical studies, for example) the dependence is explicit and clear.

## Example 2

Stratified cluster sample of individuals enrolled in the Federal Employees Health Benefit Plan (Neuhaus & Jewell (1990)).

- Response variable indicates whether someone used outpatient mental health services in the previous year for each of 1979-81.
- Clusters consist of the three values for a single person
- Four strata defined by the four possible values of the cluster totals with sampling fractions

$$f_0 = .03, \quad f_1 = \frac{1}{3}, \quad f_2 = \frac{2}{3}, \quad f_3 = 1$$

## General case

In most surveys the dependence is less clearcut with stratification determined by such things as administrative convenience or the availability of a suitable list. Often, the survey will be designed by someone else and it is not always clear why a particular stratification was chosen.

- “Even stratification could be informative to someone who knows less about the population than the survey designer” (Scott (1975))

## General case

In general, if we are not sure about the stratifying variables (or even if we are sure, but don't want to include them all in the model) we may have to fit a parametric model,  $p_h(\mathbf{y}, \mathbf{X}; \gamma)$  say, for the conditional probability of a unit being in the  $h$ th stratum given values of  $\mathbf{y}$  and  $\mathbf{X}$ .

- In our glioma example, only about one case in 70 falls in the specified sampling period.
  - natural model  $p_1(\mathbf{y}, \mathbf{X}; \gamma) = \gamma \sum_j y_j$

## General set-up

One possible strategy: include stratum indicator as a covariate in the model.

- obviously not sensible if the indicator is the response itself
- can't if we don't know the stratifying variable
- may distort the relationship of interest even if possible (eg SIDS study)

## Weighted Estimators

The standard alternative is to use weighted estimating equations.

- Stratum sizes (and hence weights) are random in most cases
  - strictly need to account for this in standard errors
- can use results from double sampling for stratification to get

$$\hat{V} \{ \hat{\boldsymbol{\theta}} \} = \left\{ \boldsymbol{\mathcal{I}}_W^{-1} + \left( \frac{1}{N} \boldsymbol{\mathcal{I}}_W \right)^{-1} \left( \sum_h W_h^2 \frac{(1-f_h)}{n_h} \hat{\boldsymbol{\Sigma}}_h \right) \left( \frac{1}{N} \boldsymbol{\mathcal{I}}_W \right)^{-1} \right\}.$$

## Weighted Estimators

- Works well if weights don't vary too much in size.
- Result always has a natural population interpretation independent of the particular sampling scheme even if the model is wrong.
- However, it can be extremely inefficient with very unequal weights (as in our example and most retrospective studies)

## Likelihood Approaches

We explore more efficient approaches based on the likelihood:

$$\begin{aligned} L(\boldsymbol{\theta}, \boldsymbol{\gamma}, g) &= \prod_{h=1}^L \left\{ \prod_{j \in D_h} \text{pr}(\mathbf{y}_{hj}, \mathbf{X}_{hj} \mid \text{unit in } \mathcal{S}_h) \right\} \text{pr}(\text{unit in } \mathcal{S}_h)^{N_h} \\ &= \prod_h \left( \prod_{D_h} \{p_h(\mathbf{y}, \mathbf{X}; \boldsymbol{\gamma}) f(\mathbf{y}_{hj} | \mathbf{X}_{hj}; \boldsymbol{\theta}) g(\mathbf{X}_{hj})\} Q_h^{N_h - n_h} \right) \end{aligned}$$

- $g(\mathbf{X})$  denotes the marginal (population) density of  $X$
- $Q_h = Q_h(\boldsymbol{\theta}, \boldsymbol{\gamma}, g)$  denotes the marginal probability that a unit is in  $\mathcal{S}_h$ .

## Likelihood Approaches

- With simple random sampling we can make inferences about  $\boldsymbol{\theta}$  conditional on the observed values of  $\mathbf{X}$  and ignore terms involving  $g(\mathbf{X})$ .
- Unfortunately, we can't ignore  $g(\mathbf{X})$  here: since
$$Q_h = Q_h(\boldsymbol{\theta}, \boldsymbol{\gamma}, g) = \int \int p_h(\mathbf{y}, \mathbf{X}; \boldsymbol{\gamma}) f(\mathbf{y} | \mathbf{X}; \boldsymbol{\theta}) g(\mathbf{X}) d\mathbf{y} d\mathbf{X}$$
  - $g(\mathbf{X})$  becomes a (potentially infinite dimensional) nuisance parameter in the likelihood.

## Likelihood Approaches

- A full likelihood approach requires us to model the covariate distribution,  $g(\mathbf{X})$ 
  - possible in simple cases (e.g. see Pearson (1903), DeMets & Halperin (1977))
  - problems involved in modelling  $g(\mathbf{X})$  are overwhelming as soon as we get past two or three covariates.
- We want to develop semi-parametric methods in which  $g(\mathbf{X})$  is left completely unspecified
  - sounds formidable - turns out to be relatively simple in practice.

## Semi-parametric approach

Let

$$\begin{aligned}\ell^*(\boldsymbol{\theta}, \boldsymbol{\gamma}, \boldsymbol{\pi}) &= \sum_h \sum_{D_h} \log(f^*(\mathbf{y}_{hj} | \mathbf{X}_{hj}; \boldsymbol{\theta}, \boldsymbol{\gamma}, \boldsymbol{\pi})) \\ &\quad - \sum_h [(N_n - n_h) \log(1 - \pi_h) + n_h \log \pi_h],\end{aligned}$$

with  $f^*(\mathbf{y}_{hj} | \mathbf{X}_{hj}; \boldsymbol{\theta}, \boldsymbol{\gamma}, \boldsymbol{\pi})$  defined by

$$f^*(\mathbf{y}_{hj} | \mathbf{X}_{hj}; \boldsymbol{\theta}, \boldsymbol{\gamma}, \boldsymbol{\pi}) \propto \pi_h p_h(\mathbf{y}_{hj}, \mathbf{X}_{hj}; \boldsymbol{\gamma}) f(\mathbf{y}_{hj} | \mathbf{X}_{hj}; \boldsymbol{\theta}).$$

The quantity  $\ell^*(\boldsymbol{\theta}, \boldsymbol{\gamma}, \hat{\boldsymbol{\pi}}_n(\boldsymbol{\theta}, \boldsymbol{\gamma}))$  is the profile likelihood obtained by maximizing the likelihood over  $\mathbf{g}(\mathbf{X})$ .

## Semi-parametric approach

Although  $\ell^*(\boldsymbol{\theta}, \boldsymbol{\gamma}, \boldsymbol{\pi})$  is not itself a likelihood, we obtain efficient semi-parametric estimators of  $\boldsymbol{\theta}$  and  $\boldsymbol{\gamma}$  by acting as if it were.

More generally, for making inferences about  $\begin{pmatrix} \boldsymbol{\theta} \\ \boldsymbol{\gamma} \end{pmatrix}$ , we can act just as if the pseudo-likelihood,  $\ell^*(\boldsymbol{\phi})$ , where  $\boldsymbol{\phi} = \begin{pmatrix} \boldsymbol{\theta} \\ \boldsymbol{\gamma} \\ \boldsymbol{\pi} \end{pmatrix}$ , is an ordinary likelihood.

## Semi-parametric approach

Specifically:

- setting the pseudo score,  $S^*(\phi) = \partial \ell^*(\phi) / \partial \phi$ , equal to zero gives the semi-parametric MLE;
- the resulting estimators are asymptotically normal;
- an estimate of  $\text{Cov}\{\hat{\theta}\}$  can be obtained directly from the appropriate block of the inverse of the observed pseudo information matrix,  $\mathcal{J}^*(\hat{\phi})^{-1}$ .

In effect, we have replaced an infinite dimensional nuisance parameter,  $g(\mathbf{X})$ , by an  $(L-1)$ -dimensional parameter,  $\pi$ . (See Alan Lee's talk tomorrow for more details.)

## Example

We have two strata here: one containing all 463 families with someone who has been diagnosed in the specified time interval, and the other containing the remaining  $N_2 = 1,942,490$  families without such a member.

- We need to fit a model for the probability of a family being in  $\mathcal{S}_1$  given  $\mathbf{y}$  and  $\mathbf{X}$ .
- Natural to model the probability of being in  $\mathcal{S}_1$  (i.e. of containing a recently diagnosed case) given  $\mathbf{y}$  and  $\mathbf{X}$  as proportional to the total number of members with brain cancer, say  $p_h(\mathbf{y}_i, \mathbf{X}_i; \gamma) = \gamma \sum_j y_{ij}$

## Results for the glioma data

	$\hat{\beta}_0(\text{se})$	$\hat{\beta}_1(\text{se})$	$\hat{\beta}_2(\text{se})$	$\hat{\sigma}(\text{se})$
Semipar:	$-5.29(.18)$	$2.12(.25)$	$0.32(.16)$	$.88(.19)$
Conditional:	—	$2.09(.36)$	$0.31(.21)$	—

## Results for the glioma data

	$\hat{\beta}_0(\text{se})$	$\hat{\beta}_1(\text{se})$	$\hat{\beta}_2(\text{se})$	$\hat{\sigma}(\text{se})$
Semipar:	$-5.29(.18)$	$2.12(.25)$	$0.32(.16)$	$.88(.19)$
Conditional:	—	$2.09(.36)$	$0.31(.21)$	—
Probands only:	—	$1.78(.39)$	$-.02(.21)$	—

## Results for the glioma data

	$\hat{\beta}_0(\text{se})$	$\hat{\beta}_1(\text{se})$	$\hat{\beta}_2(\text{se})$	$\hat{\sigma}(\text{se})$
Semipar:	$-5.60(.18)$	$2.12(.25)$	$0.32(.16)$	$.88(.19)$
Conditional:	—	$2.09(.36)$	$0.31(.21)$	—
Semip. (no N):	$-5.64(.54)$	$2.12(.26)$	$0.32(.16)$	$.00(.00)$

## Results for the glioma data

	$\hat{\beta}_0(\text{se})$	$\hat{\beta}_1(\text{se})$	$\hat{\beta}_2(\text{se})$	$\hat{\sigma}(\text{se})$
Semipar:	$-5.60(.18)$	$2.12(.25)$	$0.32(.16)$	$.88(.19)$
Semip. (no N):	$-5.64(.54)$	$2.12(.26)$	$0.32(.16)$	$.00(.00)$
Semip. ( $\gamma$ , no N):	$-5.64(.54)$	$2.12(.26)$	$0.32(.16)$	$.89(.27)$

## Results for the glioma data

	$\hat{\beta}_0(\text{se})$	$\hat{\beta}_1(\text{se})$	$\hat{\beta}_2(\text{se})$	$\hat{\sigma}(\text{se})$
Semipar:	$-5.29(.18)$	$2.12(.25)$	$0.32(.16)$	$.88(.19)$
Conditional:	—	$2.09(.36)$	$0.31(.21)$	—
Weighted:	$-5.61(.44)$	$-1.78(.55)$	$2.01(.68)$	$0(0)$

## The semi-parametric approach . . .

- can be implemented simply with a reasonably general likelihood program;
- requires absolutely no modelling of the covariate distribution;
- more efficient (sometimes much more efficient) than weighted methods;
- more robust than (and almost as efficient as?) full maximum likelihood.