

19th International Workshop on Statistical Modelling
 Florence, Italy, 4 - 8 July 2004

A Polytomous Response Multilevel Model with a Non Ignorable Selection Mechanism

Leonardo Grilli – Carla Rampichini
 Dipartimento di Statistica "G. Parenti", Firenze
 grilli@ds.unifi.it, carla@ds.unifi.it

Main goals

- Investigation of the factors which influence the way of acquisition of the professional skills;
- Comparative evaluation of the degree programmes on the basis of the adequacy of the skills they give to their graduates.

IWSM04 Grilli & Rampichini 2

Methodological issues

- The data have a **hierarchical structure**, so that the responses are correlated;
- The question of interest has a **polytomous response**;
- The response might be affected by **selection bias** due to the design of the questionnaire.

IWSM04 Grilli & Rampichini 3

Outline of the presentation

- Data description
- The model
- Main results
- Work in progress
- Conclusions

IWSM04 Grilli & Rampichini 4

The data

- Univ. of Florence, graduates of year 2000
- 3148 CATI interviews about two years after the degree
- For **employed** graduates: questions on the skills required by the present job

Do you have a job?	Frequency	Percent
1 Yes, now	2543	80.86
2 Yes, but not now	219	6.96
3 No	383	12.18
Frequency Missing = 3		

Graduates by job condition at the interview

IWSM04 Grilli & Rampichini 5

Hierarchical structure of the data

Employed graduates: distribution among the 56 degree programmes (**quite unbalanced**)

Min	Max	Med
4	386	21

Hierarchical nature of the data (graduates within course programmes)

MULTILEVEL ANALYSIS

IWSM04 Grilli & Rampichini 6

Outcome of interest

Where have you acquired the professional and technical skills?
(if used in the current job)



% Used	Way of acquisition %		
	UNIV	WORK	OTHER
91.0	43.4	36.2	11.4

The selection issue

due to the design of the questionnaire:

the outcome is missing if the skills are not used

- potential bias if the selection mechanism depends on unobserved variables correlated with the principal model's error terms
- To avoid selection bias: insert an equation that explicitly models the selection mechanism (Heckman, 1979).
- Applications rare in the multilevel framework (Borgoni e Billari, 2002).

Model with selection

(1) Selection equation Y_{ij}^S :
do you USE the skills?

(2) Acquisition equations $Y_{ij}^{P(m)}$:
WHERE have you acquired the skills?

$Y_{ij}^{P(m)}$ not observed if $Y_{ij}^S = 0$ $m=1, \dots, M$ alternatives

Not ignorable selection mechanism **IF** the two set of equations have correlated residuals

Selection equation: do you USE the skills?

Dichotomous multilevel logistic model

$$p(Y_{ij}^S = 1 | \mathbf{x}_{ij}^S, \xi_j^S, \delta_{ij}^S) = \frac{\exp(\alpha^S + \boldsymbol{\beta}^S \mathbf{x}_{ij}^S + \xi_j^S + \delta_{ij}^S)}{1 + \exp(\alpha^S + \boldsymbol{\beta}^S \mathbf{x}_{ij}^S + \xi_j^S + \delta_{ij}^S)}$$

$i=1, \dots, n_j$ graduates; $j=1, \dots, 56$ degree programmes

subjects

clusters

\mathbf{x} graduate or degree programme characteristics
 ξ random errors at degree programme level
 δ random errors at graduate level

Acquisition equations: WHERE did you acquired the skills?

Multinomial multilevel logistic model

observable only if $Y_{ij}^S = 1$

$$p(Y_{ij}^P = m | \mathbf{x}_{ij}^P, \xi_j^P, \delta_{ij}^P) = \frac{\exp(\eta_{ij}^{(m)})}{1 + \sum_{l=2}^M \exp(\eta_{ij}^{(l)})}$$

$$\eta_{ij}^{(m)} = \alpha^{P(m)} + \boldsymbol{\beta}^{P(m)} \mathbf{x}_{ij}^P + \xi_j^{P(m)} + \delta_{ij}^{P(m)}$$

$i=1, \dots, n_j$ graduates; $j=1, \dots, 56$ degree programmes
 $m=1, 2, 3$ alternatives (university, at work, other)

Nb: all parameters equal to 0 for $m=1$ (university)

\mathbf{x} graduate or degree programme characteristics
 ξ random errors at degree programme level
 δ random errors at graduate level

Assumptions on the error terms

- Errors at different levels are independent
- At each level:
 - $[\xi_j^S, \xi_j^{P(2)}, \dots, \xi_j^{P(M)}]' \sim \text{MN}(0, \Sigma_\xi)$ degree programme
 - $[\delta_{ij}^S, \delta_{ij}^{P(2)}, \dots, \delta_{ij}^{P(M)}]' \sim \text{MN}(0, \Sigma_\delta)$ graduate

If at least one of the correlations among the couples

$$(\xi_j^S, \xi_j^{P(m)}) \quad (\delta_{ij}^S, \delta_{ij}^{P(m)})$$

is not null \rightarrow not ignorable selection mechanism

IIA property

The odds for two alternatives m and l for subject i of cluster j are:

$$\frac{P(Y_{ij} = m | \mathbf{x}_{ij}, \xi_j^P, \delta_{ij}^P)}{P(Y_{ij} = l | \mathbf{x}_{ij}, \xi_j^P, \delta_{ij}^P)} = \frac{\exp[\eta_{ij}^{(m)}]}{\exp[\eta_{ij}^{(l)}]}$$

The odds depends only on the linear predictors η of the two involved alternatives and does not depend on the other alternatives.

The IIA property holds conditionally on all the covariates and error terms at all levels.

IWSM04

Grilli & Rampichini

13

Interpretation of the acquisition equations

An alternative specification of the multinomial logit model is based on the **random utility** model (McFadden, 1973):

$$U_{ij}^{(m)} = \eta_{ij}^{(m)} + \varepsilon_{ij}^{(m)}$$

$\varepsilon_{ij}^{(m)}$ are iid errors following the Gumbel distribution.

This leads to a definition of the **Intraclass Correlation Coefficient** of the m -th equation analogous to the usual definition for dichotomous logit models:

$$ICC^{(m)} = \frac{\text{Var}(\xi_j^{(m)})}{\text{Var}(\xi_j^{(m)}) + \text{Var}(\delta_{ij}^{(m)}) + \pi^2/3}$$

IWSM04

Grilli & Rampichini

14

Model identification

The parameters of the **cluster level** covariance matrix Σ_ξ are all identified;

For the **subject level** covariance matrix Σ_δ :

- the variance of δ_{ij}^S in the selection equation is obviously **not identifiable**
- the other parameters are in principle identified, but prone to empirical underidentification, unless some alternative specific covariate is included in the model (Skrondal and Rabe-Hesketh, 2003).

IWSM04

Grilli & Rampichini

15

Model likelihood

$$L(\theta | \mathbf{Y}_{ij}) = \prod_{j=1}^J \prod_{i=1}^{n_j} \left\{ \int P(\mathbf{Y}_{ij} | \mathbf{x}_{ij}, \xi_j, \delta_{ij}) f(\delta_{ij}) d\delta_{ij} f(\xi_j) d\xi_j \right\}$$

$$\theta' = (\alpha, \beta, \Sigma_\xi, \Sigma_\delta)$$

An approximated method is needed to solve the integrals

Maximization achieved by means of the **GLLAMM** command of STATA (adaptive numerical quadrature)

The subject level covariance parameters are empirically not identified (high condition number and convergence difficulties)

→ the random errors δ_{ij} are omitted from the models.

IWSM04

Grilli & Rampichini

16

Covariates selection strategy

The criterion for choosing the relevant covariates is the likelihood ratio test, with a p -value threshold of 5%.

- Selection** of the covariates **separately** for the Selection model and for the Acquisition model;
- Refinement** using the joint model, trying to reinsert in the Acquisition equations the variables which were previously discarded from the Acquisition model, but retained in the Selection model.

IWSM04

Grilli & Rampichini

17

Models comparison

	Model 1 joint S&P	Model 2 unrelated S, P
$\log L$	-2840.28	-2843.04
n. of parameters	26	24
<i>Random parameters</i>		
$\text{Var}(\xi_j^S)$	0.1510	0.1448
$\text{Var}(\xi_j^{P(2)})$	0.1637	0.1531
$\text{Var}(\xi_j^{P(3)})$	0.4274	0.4129
$\text{Corr}(\xi_j^S, \xi_j^{P(2)})$	-0.2401	.
$\text{Corr}(\xi_j^S, \xi_j^{P(3)})$	-0.6772	.
$\text{Corr}(\xi_j^{P(2)}, \xi_j^{P(3)})$	0.8768	0.8479

LR Test
 $\chi^2=5.52$, $df=2$
 $p\text{-value}=0.0633$



without selection model

Estimation with
GLLAMM

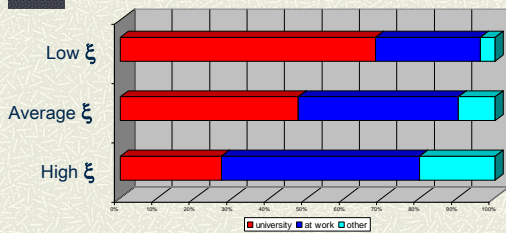
Fixed parameters quite similar in Models 1 and 2!

IWSM04

Grilli & Rampichini

18

Probabilities of acquisition by degree programme (baseline graduate)



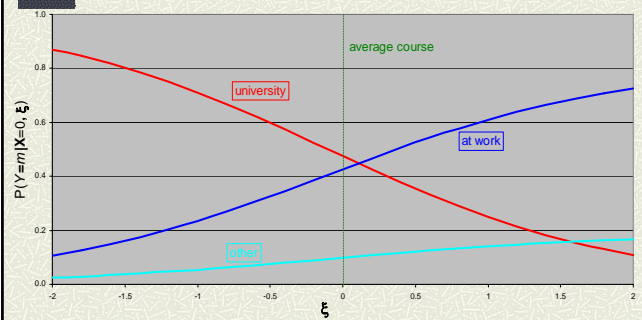
Given graduate and degree programme observed characteristics, the probability of each alternative vary a lot among degree programmes!

IWSM04

Grilli & Rampichini

19

Acquisition probabilities and ξ values ($\xi^{P(2)} = \xi^{P(3)}$, baseline graduate)

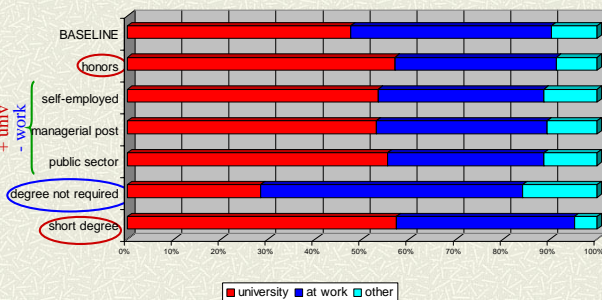


IWSM04

Grilli & Rampichini

20

Acquisition probabilities by graduate characteristics (average degree prog.)

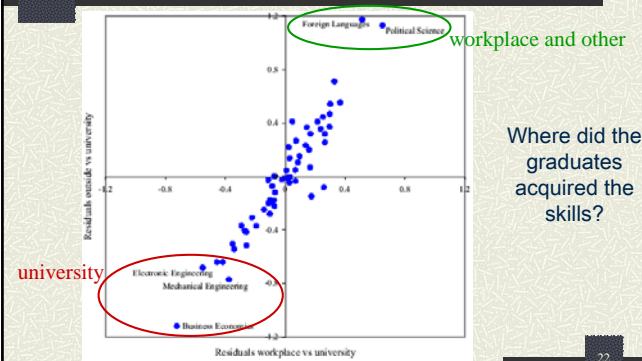


IWSM04

Grilli & Rampichini

21

Degree programme level EB residuals from the acquisition equations



22

Work in progress...



- Further work is needed to fully understand the implications of selection mechanisms that act in a hierarchical framework;
- Some sensitivity analysis is required to assess how the results are affected by the distributional assumptions.

IWSM04

Grilli & Rampichini

23

Work in progress (cont'd)



Simulation study:

- How do the selection mechanism interact with the hierarchical structure of the data?
- What are the consequences of ignoring the hierarchy?
- How change the power of the LR Test for presence of selection?
- What is the influence of the distributional assumptions on error terms?

IWSM04

Grilli & Rampichini

24

Preliminary results ...

We start considering a two-level bivariate probit model with selection mechanism and covariates at both levels (some are common).



Given a certain total correlation between the two latent variables, **things are better if the correlation is mainly due to the clusters.**

Estimation

- ✦ The ML estimation algorithm based on **adaptive numerical quadrature**, used in the application, is accurate and flexible, but it requires long computational times, which rapidly increase with the model complexity.
- ✦ Many alternative estimation methods are possible, e.g. Bayesian MCMC and Maximum Simulated Likelihood (Train, 2003).

Concluding remarks

- ✦ We have shown how to build a complex polytomous response model for the analysis of graduates' skills, taking into account:

- the hierarchical structure of the phenomenon
- the adjustment for a possible selection bias.

In the application the hierarchical structure has a crucial role, while selection bias results negligible.

- ✦ The model allows to:

- characterize ways of acquisition of the skills
- find out extreme degree programmes to be further investigated.

Thanks for your attention...



References

- Heckman J.J. (1979) Sample selection bias as a specification error, *Econometrica*, 47, 153–161.
- McFadden D. (1973) Conditional logit analysis of qualitative choice behaviour, in: *Frontiers in Econometrics*, Academic Press, New York.
- Rabe-Hesketh S., Pickles A. and Skrondal A. (2001) GLLAMM manual, Technical Report 2001/01, Department of Biostatistics and Computing, Institute of Psychiatry, King's College, London.
- Skrondal A. and Rabe-Hesketh S. (2003) Multilevel logistic regression for polytomous data and rankings, *Psychometrika*, 68, 267–287.
- Train K. (2003) *Discrete Choice Methods with Simulation*, Cambridge University Press, New York.