Estimating linear and generalized linear models using Respondent Driven Sampling (RDS)

Märt Möls¹,

Krista Fischer², Anneli Uusküla¹ ¹ Tartu University, Estonia; ² MRC Biostatistics Unit

Estonian Intravenous Drug Users Survey (Data collection 2005)

In co-operation with Imperial College London (Lucy Platt, Natalia Bobrova and Tim Rhodes)

- Estonian Commercial Sex Workers Survey (Data collection 2006)
- Estonian Intravenous Drug Users (IDU) Survey 2 (Data collection 2007)
 In co-operation with Don DesJarlais (Beth Israel Medical Center), partly funded by CRDF and NIH (NIDA).

At least 125 studies worldwide using RDS methodology – drug users, men having sex with men, sex workers, high risk heterosexual men, homeless people, jazz musicians,



















Basic description of the Studies

	Intervenous Drug Users	Commercial Sex Workers	IDU 2
Sample size	450 Tallinn - 350; Kohtla-Järve-100	227	700
Number of seeds	9	43	11
Median network size	50	6	50
Number of "waves"	9	9	17
HIV+	62%	7,5%	63%

Wishlist from epidemiologists

- Prevalence of HIV+ (probability to be HIV+), mean duration of injection carieer, ... (with Confidence Intervals)
- Test wether sharing needles/type of drug/ ... is a risk factor for HIV+; how much the odds to get HIV increase if someone adopts riskier behaviour (tests, Odds Ratios OR with CI)
 Models for predicting the probability of HIV;

How to analyse a RDS study?

- Naive Analysis Treat the data as from usual random sample. Convinient, often used, but...;
- Heckathorn's method and software use Markov Chains to derive asymptotically unbiased (under plausible assumptions) estimates of proportions;
- Use Linear Mixed Models/Generalized Linear Mixed Models with suitable covariance structure (our proposal).

Naive Analysis - problems

- Undersampling of respondents with few friends (small network size);
- 2. Bias due to non-random selection of seeds;
- 3. Friends are "more similar" to each other than two randomly choosen persons from the target population

Association between recruiter and recruited (from IDU study):

HIV status	p-value=0,0009
age	p-value=0,03
type of drug used	p-value=0,00002

- "Respondent-Driven Sampling: A New Approach to the Study of Hidden Populations." By Douglas D. Heckathorn. Social Problems, 1997.
- "Respondent-Driven Sampling II: Deriving Valid Population Estimates from Chain-Referral Samples of Hidden Populations." By Douglas D. Heckathorn. Social Problems, 2002.
- www.respondentdrivensampling.org

Nearly-unbiased estimation of a distribution with few possible categories (for example: infected/not-infected), under understandable set of assumptions.

Bootstraping methods to calculate approximate standard errors and confidence intervals for prevalence estimates

Heckathorn's method: limitations

Limited possibilities for more complex models; Limited possibilities to handle continous variables;

Now we will continue with our proposals (which are actually just some old classical methods) There is nothing new in analysing correlated observations.

For example, one can use

Time Series Analysis



Repeated Measures / Multilevel Analysis



Analysing Correlated Data

- Estimate the correlation/covariance matrix V by using ML or REML;
- Estimate parameter vector via Estimated BLUE (GLS):

$$\hat{\boldsymbol{\beta}} = (\boldsymbol{X}^T \hat{\boldsymbol{V}}^{-1} \boldsymbol{X})^{-1} \boldsymbol{X}^T \hat{\boldsymbol{V}}^{-1} \boldsymbol{Y}$$

Test hypothesis etc:

$$L\hat{\beta} \pm t_{\hat{\nu},\alpha/2} \sqrt{L\hat{C}L^T}$$

$$F = \frac{\hat{\beta}^T L^T (L\hat{C}L^T)^{-1} L\hat{\beta}}{\operatorname{rank}(L)}$$

$$\hat{C} = (X^T \hat{V}^{-1} X)^-$$

See:

SAS Online Documentation, PROC MIXED, PROC GLIMMIX;

Searle et al – Variance Components; Generalized, Linear, and Mixed Models

Other textbooks about Mixed models/ linear models

A possible correlation structure for a RDS Design



	Α	B	С	D
Α	1	r	r^2	r ³
В	r	1	r	r^2
C	r^2	r	1	r
D	r^3	r^2	r	1

Possible correlation structure for a RDS Design



	Α	B	С	D	E	F	G	Η	Ι
Α	1	r	r^2	r ³	r	r^2	r ³	r ³	r ³
В	r	1	r	r^2	s	r	r^2	r^2	r^2
С	r^2	r	1	r	Sr	s	Sr	SV	Sr
D	r ³	r^2	r	1	sr ²	SV	sr ²	sr ²	sr ²
Ε	r	S	Sr	sr ²	1	SV	sr ²	sr ²	sr ²
F	r^2	r	s	Sr	Sr	1	r	r	r
G	r^3	r^2	Sr	sr ²	sr ²	r	1	s	s
Η	r^3	r^2	SV	sr ²	sr ²	r	S	1	S
Ι	r^3	r^2	Sr	sr ²	sr ²	r	S	S	1

Possible correlation structure for a RDS Design



Logistic regression / Generalized Linear Mixed Model SAS PROC GLIMMIX documentation: (TW = Y), TV = c1(Y)

 $g(EY) = X\beta; \qquad EY = g^{-1}(X\beta)$ $Var(Y) = A^{1/2}VA^{1/2}$

The matrix A is a diagonal matrix and contains the variance functions of the model. The variance function expresses the variance of a response as a function of the mean.

Undersampling

1. Undersampling can be corrected using weighted averages or related techniques (Horwitz-Thompson estimator). Corrections are regularly applied if stratified random sampling has been used, for example.

 $P(Sampled \mid Network \quad Size = i) \propto i$

Correction for undersampling

- 1. Include network size to the model of interest, eg.: logit(P(HIV+|NS)) = c+f(NS)
- 2. Integrate *NS* out from the final result based on the estimated proportions of network sizes, eg: $P_{est}(HIV+) = \sum_{i} P_{est} (HIV+|NS=i) P_{est} (NS=i)$
- 3. One can use for example delta method to calculate standard error for the estimate (is some better method available?)

Correction for undersampling

A note:

Most applications in epidemiology (probably) do not require the complicated procedure to calculate the standard error. If one is interested in estimating the effect of a risk factor and there is no interaction effect between the risk factor (RF) and network size, one can just fit a model

logit(P(HIV+|NS)) = c+f(NS)+RF+confounders

to the data (+correct covariance structure) and the delivered inference for RF is still valid.





Seeds are correlated, friends are correlated, but correlation goes to zero as the distance within recruitment chain increases.

Seed selection II



Different non-communicating populations, one recruitment chain will never jump to a different population (respondents within one recruitment chain will always remain correlated), but seeds can be considered as a random sample from target population

Seed selection bias III – Problem!



Different non-communicating populations, one recruitment chain will never jump to a different population, seeds are also correlated – not truely a random sample from target population

Does it work?

Comparison with other methods
Simulations
Does the proposed model fit to the real data?

Some results I

CSW Study:

	Naive	Heckathorn	GLMM
HIV+ $(\%)$	7,5%	4,7%	4,7%
se	0,01747	_	0,0169
95%-CI	4%12%	3%9%	1%8%

Some results II

IDU Study:

	Naive	Heckathorn	GLMM
Tallinn, HIV+ (%)	54%	47%	47%
Kohtla-Järve, HIV+ (%)	89%	91%	90%
Average Age	24,2	NA	23,7

Simulation results (True population prevalence about 0,45)

95%-CI

Method	MSE	coverage	bias
Naive	0,313	76,0%	0,046
Heckatorn	0,141	96,5%	0,003
GLMM	0,143	96,0%	-0,006

Does the model fit?

Example: linear model (for age) from the IDU study

Correlation parameter estimates: r = 0,084s = 0,152

AIC = 2700(RDS correlation structure)AIC = 2704(Independence)

Conclusions

- Regression and logistic regression models can be fit and valid inference can be made (for RDS samples)
- One can correctly analyse an RDS study by using standard software (for example R)

 However: untestable(?) assumptions; using the software can be tricky (lack of documentation) etc.

Few additional slides, just in case...

Estonian IDU Study



HIV risk factors - CSW

	Naive Analysis	GLMM
	OR (95% CI)	OR (95% CI)
	(unadjusted)	(unadjusted)
Category of SW		
Brothel	4,0 (1,215,3)	3,6 (1,210,4)
Street	12,6 (3,056,8)	3,1 (0,714,1)
Other	1	1
Drug use: No	1	1
Yes	8,3 (2,327,8)	2,3 (0,68,5)
Years of CSW	0,84 (0,700,95)	0,76 (0,670,87)

HIV risk factors - IDU

	Naive Analysis	GLMM
	OR (95% CI)	OR (95% CI)
	(unadjusted)	(unadjusted)
Gender: male	1	1
female	1.2 (0.72.1)	1.3 (0.82.3)
Years of drug use:		
2 years or less	1	1
35 years	2.8 (1.45.5)	2.4 (1.24.9)
69 years	4.8 (2.49.3)	3.8 (1.97.4)
10 years or more	3.5 (1.77.0)	2.5 (1.25.2)

IDU – final model for HIV status

	log(OR) (DR 95%-CI	p-value
Network size (Ref: 100+)			
less than 100	-0.75 0.	47 0.300.7	75 0.002
Duration of injection career	(Ref: 0-2 year	rs)	
3-5 years	0.94 2.	55 1.165.62	2 0.020
6-9 years	1.61 5.	00 2.1911.38	3 0.001
10 years or more	1.31 3.	69 1.429.61	0.008
Place of residence (Ref: Talli	nn)		
Kohtla-Järve	1.79 5.	99 2.5314.19	0.005
Number of sexual partners d	uring last yea	ar (Ref: 0)	
one	-1.37 0.	25 0.090.69	0.007
more than one	-1.07 0.	34 0.130.89	0.028
Age group (Ref: <20)			
20-24	-0.43 0.	65 0.331.29	0.218
25-29	-0.90_0.	41 0.180.89	0.026
30 or more	-1.26_0.	28 0.110.75	5 0.012

Software

- Most major statistical software packages (SAS, STATA, S-plus, R,...) can correctly analyse data with correlated errors – one just needs to specify correct correlation structure.
- However, no package offers a ready-made correlation structure suitable for analysing RDSdesigns. Some packages allow user to supply their own user-specified correlation structure (R, S-Plus)

Software example: R

R is a free statistical package (www.r-project.org)

- To create a new correlation structure (corNew) one has to write five new functions:
 - 1. corNew constructor function creates a new object
 - 2. Initialize.corNew initializes the new object
 - 3. corMatrix.corNew returns the correlation matrix (based on current parameter values)

4. Functions "coef.corNew" and "coef<-.corNew" to extract and change the correlation parameters.

Created correlation structure can then be used in other functions (gls, lme, glmmPQL,...)



2. Estimate the average network size for all groups



3. Assumption:

The total number of ties from HIV- group to HIV+ group is the same as the total number of ties from HIV+ group to HIV- group (If A knows B then B also knows A).

 $N_+ ANS_+ p_{+-} = N_- ANS_- p_{-+}$

From where:

 $N_{+}/(N_{+}+N_{-}) = \frac{ANS_{-}p_{-+}}{ANS_{-}p_{-+}+ANS_{+}p_{+-}}$

 $3,40 \cdot 0,054 / (5,12 \cdot 0,769 + 3,40 \cdot 0,054) = 0,045$

Distribution of network sizes

$$P(S = 1 | NS = i) = \frac{P(NS = i | S = 1)P(S = 1)}{P(NS = i)}$$

$$P(NS = i \mid S = 1) \propto P(NS = i) i$$

Problems with the CSW-study



Problems



o		ο		0		0
	0		0		0	
0		0		0		0
	0		0		0	
0		0		0		0
	0		0		0	
0		0		0		0















