

**FUNDAÇÃO GETULIO VARGAS  
SCHOOL OF APPLIED MATHEMATICS**

**LUCAS MACHADO MOSCHEN**

**PREVALENCE ESTIMATION AND BINARY REGRESSION  
METHODS FOR RESPONDENT-DRIVEN SAMPLING WITH  
OUTCOME UNCERTAINTY**

Rio de Janeiro

2021

# Contents

1	INTRODUCTION . . . . .	4
2	THEORETICAL BACKGROUND . . . . .	6
2.1	Prevalence estimation problem . . . . .	6
2.2	Respondent-driven sampling . . . . .	10
2.2.1	Details about the sampling procedure . . . . .	11
2.2.2	Assumptions and statistical properties . . . . .	13
2.2.3	Models for the RDS Process . . . . .	14
2.2.3.1	First-order Markov process . . . . .	15
2.2.3.2	Successive sampling (SS) . . . . .	16
2.2.3.3	Graphical Structure model . . . . .	16
2.2.3.4	New model . . . . .	19
2.2.4	Prevalence estimators . . . . .	19
2.2.5	Regression methods . . . . .	20
2.2.6	Bootstrap methods for uncertainty quantification . . . . .	21
2.2.7	Diagnosis of RDS . . . . .	21
2.3	Generalized linear models . . . . .	21
2.4	Conditionally autoregressive models . . . . .	22
2.5	Bayesian statistics . . . . .	22
2.6	Computational methods . . . . .	23
2.6.1	Hamiltonian Monte Carlo . . . . .	23
2.6.1.1	Diagnostics . . . . .	24
2.6.2	Metropolis-within-Gibbs . . . . .	24
3	PREVALENCE MODELLING AND REGRESSION METHODS . . . . .	25
3.1	Perfect tests . . . . .	26
3.1.1	Identifiability . . . . .	26
3.1.2	Simulated data . . . . .	29
3.2	Sensitivity and specificity . . . . .	32
3.2.1	Independent beta distribution priors . . . . .	33
3.2.2	Hierarchical partial pooling prior . . . . .	33
3.2.3	Bivariate Beta prior . . . . .	34
3.3	Imperfect tests . . . . .	35

3.3.1	Simulated data . . . . .	36
3.4	Imperfect tests and respondent-driven sampling . . . . .	36
3.4.1	Simulated data . . . . .	37
3.4.2	Exponential Random Graph Model (ERGM) . . . . .	37
3.5	Model extensions . . . . .	38
3.6	Mispecified data simulation . . . . .	39
4	<b>DISCUSSION ABOUT PRIOR DISTRIBUTIONS AND SENSITIVITY ANALYSIS</b> . . . . .	40
4.1	Prior analysis of sensitivity and specificity . . . . .	40
4.2	Prior analysis on the parameter tau . . . . .	40
4.3	Prior analysis on theta . . . . .	40
5	<b>REAL DATA APPLICATIONS</b> . . . . .	41
6	<b>CONCLUSION</b> . . . . .	42
	<b>References</b> . . . . .	43
	<b>APPENDIX</b> . . . . .	49
	<b>APPENDIX A – BIVARIATE BETA DISTRIBUTION</b> . . . . .	50
A.1	Comments about integration . . . . .	53
A.2	Specifying parameters $\alpha$ . . . . .	54
	<b>APPENDIX B – STAN CODES</b> . . . . .	58

# Todo list

Fix order after.	6
Should I mention more reasons to study prevalence?	6
It might be nice to add examples.	6
Provide some reference	7
It may be good to justify this choice.	7
Include notation of RDS used posteriorly.	14
Provide an example to explain all the above definitions.	17
It would be nice to cite Camila, but reference not found.	21

## **List of sections to revise**

1. Respondent-driven sampling;
2. Add Hierarchical modelling chapter;
3. Should I add a subsection in Bayesian Statistics revising Prevalence estimation models using Bayesian paradigm?

## **What to do after?**

1. Notes about Bivariate Beta;
2. Study case about CAR models in bernoulli aspect.

# 1 Introduction

Hidden or hard-to-reach populations have two main features: no sampling frame exists, given that their size and boundaries are unknown, and there are privacy concerns because the subjects are stigmatized or have illegal behavior (HECKATHORN, 1997). Fear of exposition or prosecution complicates the enumeration of the populations and the learning about them. Moreover, if the occurrence frequency of the condition is low, there are high logistic costs involved. Some examples are heavy drug users, sex workers, homeless people, and men who have sex with men.

Research has been carried out with the development of some methods to reach these populations, such as, for example, snowball sampling (GOODMAN, 1961), key important sampling (DEAUX; CALLAGHAN, 1985), and targeted sampling (WATTERS; BIERNACKI, 1989). (HECKATHORN) introduced the Respondent-Driven Sampling (RDS) to fill some gaps from other methods he depicted in his work. In his proposed approach, the researchers select a handful of individuals from the target population and give them coupons to recruit their peers. The individuals receive a reward for being recruited and for recruiting, which creates a dual incentive system. After (HECKATHORN, 1997), several papers studied this topic more deeply.

Following the sampling from the target population, a questionnaire or a disease test is conducted. This work considers binary outcomes. For instance, asking about smoking status or testing for HIV infections. However, the diagnoses are subject to measure error, and regard their accuracy is a vital step (REITSMA et al., 2005). One common way to do this is to measure jointly *sensitivity* and *specificity*. The former is the ability to detect the condition, while the latter to identify the absence of it.

Nevertheless, because of our lack of knowledge about Nature itself, it is necessary to model the uncertainty of this process, and Bayesian Statistics is the indicated area of study. In the Bayesian paradigm, the parameters are random variables, and the beliefs about them are updated given new data. The idea is to propagate uncertainty about the outcome through the network of contacts, which has its probability distribution.

This work proposes to study the survey method Respondent-Driven Sampling (RDS), a chain-referral method with the objective of sampling from hard-to-reach

populations when necessary to estimate the prevalence of some binary condition from this population. The modeling also accounts for sensibility and sensitivity since the imperfection of the detection tests. We also intend to apply this framework efficiently, comparing Monte Carlo algorithms and Laplace approximations.

## 2 Theoretical background

In this chapter, we shall describe the theoretical background taken under consideration for the developed models and analysis, including Bayesian statistics (Section 2.5), the prevalence estimation problem (Section 2.1), Respondent-driven sampling (Section 2.2), and computational methods (Section 2.6) used in our research.

Fix order after.

### 2.1 Prevalence estimation problem

The study of how health-related conditions are distributed among populations is known as *Epidemiology* (ROTHMAN; GREENLAND; LASH, et al., 2008, p. 32), which aims to derive valid estimates for potential causes from diseases that affect people. It is a fundamental research area in policy formulation, implementation of prevention programs, and development of laws. In order to accomplish these goals, the epidemiologists use some *measures of disease frequency*, including *incidence* and *prevalence*. The former is related to the proportion of new cases of a disease given a period of time, while the latter is the proportion of individuals exposed at time  $t$  and it is the object of study of this section. An interesting point is the following:

Diseases with high incidence rates may have low prevalence if they are rapidly fatal or quickly cured. Conversely, diseases with very low incidence rates may have substantial prevalence if they are nonfatal but incurable. (ROTHMAN; GREENLAND; LASH, et al., 2008, p. 46).

As a result, prevalence represents both incidence and the duration of disease. Noordzij et al. (2010, p. c18) highlights that prevalence reveals the burden of a disease in respect to its effects on society, such as, monetary costs, quality of live, and morbidity. They also comment that when measured periodically, its evolution can identify potential causes of the infection and prevention and care methods. We remark that when it is impossible to test every individual at the same time, we assume that all individuals remain exposed to the disease at time of the last tested individual.

Should I mention more reasons to study prevalence?

Consider a population of interest and a known condition, such as, for instance, a disease or a binary behavior. A diagnostic test is done in the individuals to measure the presence or the absence of this condition, such as serological tests.

It might be nice to add examples.

Mathematically, we denote  $\theta \in (0, 1)$  the prevalence of the condition, which is the parameter of interest. Let  $I$  be a index set for the individuals. We also denote  $Y_i^{\text{true}}$  the indicator function of the presence of the condition in the  $i^{\text{th}}$  individual, that is,

$$Y_i^{\text{true}} = \begin{cases} 1, & \text{if individual } i \text{ has the condition.} \\ 0, & \text{otherwise.} \end{cases}$$

Assume for simplicity that all tests are performed at time  $t$ . Assume that  $Y_i$  indicates the result of the test, then

$$Y_i = \begin{cases} 1, & \text{if test was positive in individual } i. \\ 0, & \text{otherwise.} \end{cases}$$

Since it is not usually feasible to test everyone in the population, it is necessary to random select individuals from the population. On that point, other sampling approaches may be better options, such as stratified random sampling, systematic sampling, and two-stage cluster sampling. From that experiment, we get a sample  $y = \{y_1, \dots, y_n\}$ . Based on that outcomes the Maximum Likelihood Estimator is the following expression

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n y_i, \quad (2.1)$$

which is an estimator for the *apparent prevalence*, that is, the probability of a positive outcome.

However, this estimator assumes that the diagnostic test used is perfect, which is often incorrect. It is also not interesting when the samples are not randomly selected (See Section 2.2). From that point, it is crucial to regard the evaluation of the diagnostic procedure by some measurement. Šimundić (2009, p. 2) presents several options with different aspects, such as the *likelihood ratio*, *sensitivity and specificity*, and *the area under the ROC curve*. In this work, we consider the sensitivity and specificity of the test.

Provide some reference

A perfect test would discriminate every sick individual from the non-sick ones. Given that there is not such thing, we suppose having a *gold standard test* that is the best available test (VERSI, 1992) to diagnose a particular disease. Its result is a proxy for the real  $Y_i^{\text{true}}$  and

It may be good to justify this choice.

In the context of infectious diseases, a gold standard can be a very precise molecular test that detects the presence of the pathogen's genetic material, polymerase chain reaction (PCR) for instance. (BASTOS; CARVALHO; GOMES, 2021, p. 125).

From the gold standard, we can evaluate a second test, typically faster or cheaper. The possible results upon comparing these tests are presented in table 1. The definitions for each initials in the table are the following:

- a) true positive (TP): when both tests agree that the individual has the disease;
- b) true negative (TN): when both tests agree that the individual does not have the disease;
- c) false positive (FP): when the test under evaluation has a positive diagnose, despite the golden standard being negative;
- d) false negative (FN): when the test under evaluation has a negative diagnose, despite the golden standard being positive.

Chart 1 – Two-by-two table that compares the result from the gold standard to the test under evaluation.

	$Y = 0$	$Y = 1$
$Y^{\text{true}} = 0$	TN	FP
$Y^{\text{true}} = 1$	FN	TP

Source: Prepared by the author (2021) and based on [Bastos, Carvalho, and Gomes \(2021, p. 126\)](#).

*Remark 2.1.1.* When a gold standard test is not available, which is called no gold standard situations ([RUTJES et al., 2007, p. 1](#)), other methods should be considered such as the construction of reference standard by giving the patients either different or the same tests and combining the results somehow. For more details, [Rutjes et al. \(2007\)](#) does a literature review on the topic.

For now, we drop the index  $i$  in the random variables  $Y_i$  and  $Y_i^{\text{true}}$ . Let  $p = \Pr(Y = 1)$  be the probability of a positive test. We call  $p$  the *apparent prevalence* since it is what the researchers observe. Equation (2.1) is an estimator for it. We also have that  $\Pr(Y^{\text{true}} = 1) = \theta$ . Notice that  $p$  depends on the used test, while  $\theta$  does not. In prevalence estimates, we will only have  $\theta = p$  if the test is perfect or the test is the gold standard itself. Define the following:

**Definition 2.1.1** (Sensitivity). Probability of a positive test correctly identified. In mathematical terms, conditioned on  $Y^{\text{true}} = 1$ , the *sensitivity*  $\gamma_s$  is the probability of  $Y = 1$ :

$$\gamma_s = \Pr(Y = 1 | Y^{\text{true}} = 1). \quad (2.2)$$

**Definition 2.1.2** (Specificity). Probability of a negative test correctly identified. In mathematical terms, conditioned on  $Y^{\text{true}} = 0$ , the *specificity*  $\gamma_e$  is the probability of

$Y = 0$ :

$$\gamma_e = \Pr(Y = 0 | Y^{\text{true}} = 0). \quad (2.3)$$

**Theorem 2.1.1** (Relation between prevalence and apparent prevalence). *These quantities are related by the following equation:*

$$p = \gamma_s \theta + (1 - \gamma_e)(1 - \theta). \quad (2.4)$$

*Proof.* This is a direct application of the definition of conditional probability and the countable additivity axiom of Probability:

$$\begin{aligned} p &= \Pr(Y = 1) = \Pr(Y = 1, Y^{\text{true}} = 1) + \Pr(Y = 1, Y^{\text{true}} = 0) \\ &= \Pr(Y = 1 | Y^{\text{true}} = 1) \Pr(Y^{\text{true}} = 1) + \Pr(Y = 1 | Y^{\text{true}} = 0) \Pr(Y^{\text{true}} = 0) \\ &= \Pr(Y = 1 | Y^{\text{true}} = 1) \Pr(Y^{\text{true}} = 1) \\ &\quad + (1 - \Pr(Y = 0 | Y^{\text{true}} = 0))(1 - \Pr(Y^{\text{true}} = 1)) \\ &= \gamma_s \theta + (1 - \gamma_e)(1 - \theta). \end{aligned}$$

□

The intuition behind this equation is pretty simple: the proportion of positive test counts the correct identified exposed individuals and the incorrect identified not exposed. Equation (2.4) also reveals that if  $\gamma_s = \gamma_e = 1$ , we have the trivial case  $p = \theta$ . Moreover, if  $\gamma_s = \gamma_e = 0.5$ , we have that  $p = 0.5$  and there is no information about  $\theta$ .

A frequentist approach assumes that  $\theta$  is fixed and unknown. Its inference is based on the point estimate for the apparent prevalence  $\hat{p}$  given in Equation (2.1), along with a Confidence Interval, such as the Wald Confidence Interval built with a normal approximation. In order to provide a point estimate for  $\hat{\theta}$ , Rogan and Gladen (1978, p. 73) propose

$$\hat{\theta}^{RG} = \frac{\hat{p} - (1 - \gamma_e)}{\gamma_s + \gamma_e - 1}. \quad (2.5)$$

Suppose a disease with prevalence  $\theta = 0.01$ . In this case, we would have that  $p \approx 1 - \gamma_e$  by equation (2.4). Given the randomness, it is possible to have  $\hat{p} < 1 - \gamma_e$ , which would define a useless estimative for  $\theta$ . Besides that, Confidence Intervals for that expression does not include uncertainty about  $\gamma_e$  and  $\gamma_s$ . On the other side, a Bayesian approach let  $\theta$  be a random variable, allowing the researcher to incorporate their uncertainty on the prior distribution, which is explained in Section 2.5. It also allows to include uncertainty in sensitivity and specificity of the test. According to Branscum, Gardner, and Johnson (2005):

Diagnostic-test evaluation is particularly suited to the Bayesian framework because prior scientific information about the sensitivities and specificities of the tests and prior information about the prevalences of the sampled populations can be incorporated. (BRANSCUM; GARDNER; JOHNSON, 2005, p. 1).

Therefore, this work focus on the Bayesian paradigm.

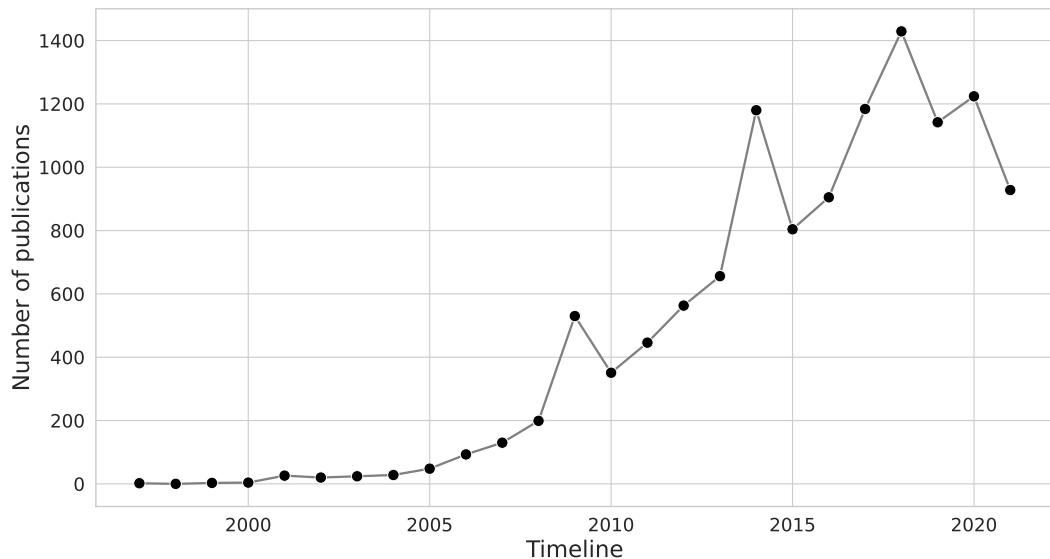
## 2.2 Respondent-driven sampling

Respondent-driven sampling (RDS) is a procedure developed by Heckathorn (HECKATHORN, 1997) to survey *hidden* or *hard-to-reach populations*, whose main characteristic is the absence of a sampling frame, i.e., it is not possible to enumerate its individuals since size and boundaries are unknown. The second characteristic of these populations is the confidentiality concerns, given that membership is stigmatized or illegal. With that aspect, traditional sampling methods which produce probability samples are infeasible. To overcome this, Snowball Sampling (GOODMAN, 1961) is the most common method, and it relies on the respondents to nominate more subjects within the population as a snowball. Examples of studied groups include people who inject drugs (PWID), men who have sex with men (MSM), and female sex workers (FSW) (GILE; BEAUDRY, et al., 2018, p. 66).

Heckathorn's proposal (1997) was to specialize this method without the need of nominating peers. In this approach, the researchers select some individuals, called *seeds* from the target population, and give them a fixed amount of *recruitment coupons* to recruit their peers. Each recipient of the coupons reclaims it in the study site, is interviewed, and receives more coupons to continue the recruitment. This process occurs until it reaches some stopping criteria, such as the sample size achieving some desired number. The sampling is without replacement, so the participants cannot be recruited more than once. Moreover, the respondents inform how many subjects from the population they know. Other less usual methods include Key Important Sampling (DEAUX; CALLAGHAN, 1985), and Targeted Sampling (WATTERS; BIERNACKI, 1989), both are convenience sampling methods.

According to Gile, Beaudry, et al. (2018, p. 66), there are two main advantages of RDS over other snowball samplings. First, the fixed number of recruitment coupons enforces the network gets deeper and distant from the seeds, which reduces the dependence of the final sample from the initial chosen by researchers. Second, since the recruited subjects do not have to name their peers, confidentiality is maintained until the recruitment is completed. Other problems cited by Heckathorn (1997, p.

Figure 1 – Publications by year with the term “Respondent driven sampling” from 1997 to 2021.



Source: <https://app.dimensions.ai>. Exported on October 31, 2021.

175) include biases towards individuals who are more cooperative, biases by masking when the participants do not name friends for the next wave to protect them, and individuals with more links may be oversampled. RDS offers a solution with a *dual incentive system*, explained in Subsection 2.2.1.

Since the creation of the method by Heckathorn, several papers have been published, as Figure 1 presents. The figure was produced searching publications with the term “Respondent-driven sampling.” These works generally aim to give basis to public health policies. Good examples in Brazil are (DAMACENA et al., 2019), (MOTA, 2012), and (BASTOS; BASTOS, et al., 2018). Damacena et al. (2019) apply the RDS method to carry out biological and behavioral surveillance in FSW populations from twelve cities in Brazil. Mota (2012) proposes the RDS method in MSM populations from ten cities in Brazil. Bastos, Bastos, et al. (2018) study several sexually transmitted infections among transgender women from twelve Brazilian cities.

### 2.2.1 Details about the sampling procedure

The RDS method was expanded by Heckathorn (2002). It detailed two aspects: introducing a way to correct *homophily* biases that is the tendency for individuals to connect to others similar to them, and *personal network size* or *degree* that is the number of connections of an individual within the target population. It also presented

a bootstrapping procedure to quantify uncertainty about inferences. [Salganik and Heckathorn \(2004\)](#) slightly modified the RDS procedure and introduced proof that under some regularity conditions, RDS estimators were asymptotically unbiased. [World Health Organization \(2013\)](#) is a reference to know how to execute an RDS survey. According to it:

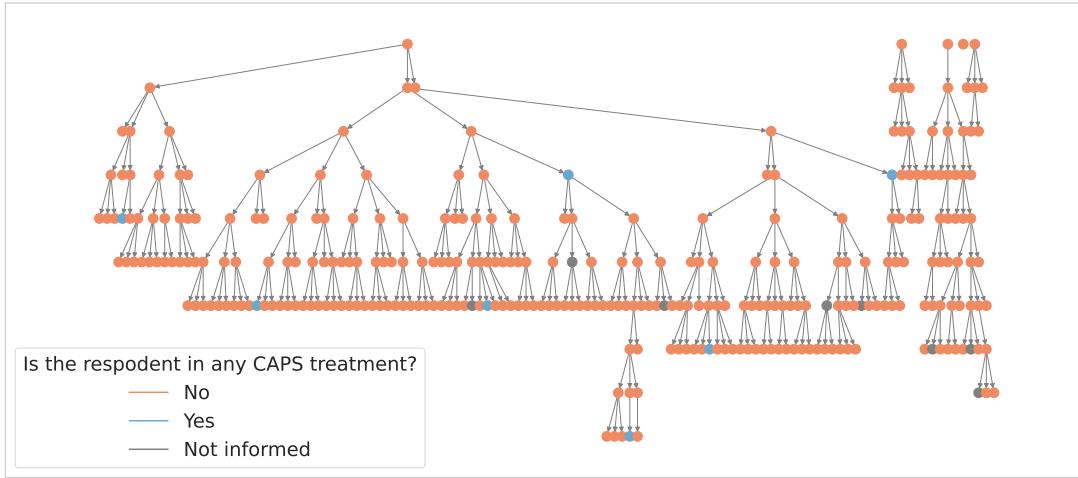
Seeds are non-randomly selected members of the survey population who initiate the RDS recruitment process. From each seed, a recruitment chain is expected to grow. Seeds play an extremely important role in conducting an RDS survey. ([WORLD HEALTH ORGANIZATION, 2013](#), p. 70).

No rule was established on the number of seeds to start the sampling. It typically varies from 2 to 32, with the mean being 10 ([WORLD HEALTH ORGANIZATION, 2013](#), p. 70). The number can not be small since unsuccessful recruitments are common. A diverse choice among the target population may accelerate the convergence to equilibrium. It also allows the access to isolate and subpopulations. After this selection, three coupons are distributed to each participant. The coupons must have information about survey site location, a unique identification code, telephone number, and opening hours. [Gile, Beaudry, et al. \(2018](#), p. 67) highlights that “this number is chosen to strike a balance between the inferential desire [...] and the practical necessity of guarding against early termination of the sample trees.”

Subjects receive a reward for being interviewed and recruiting their peers within the target population, which establishes a dual incentive system. The *primary incentive* is the *individual-sanction-based control*, so there is a reward for participating in the survey. The second one is the *group-mediated social control* that influences the participants to induce others to comply to get the remuneration for the recruitment. When social approval is relevant for the members, recruitment can be more efficient and cheaper. It happens because material incentives are converted into peer-based symbolic since there is social influence involved. In conclusion, consenting to be recruited provide material and symbolic motivation to both recruiter and participant.

For an illustrative example, [Figure 2](#) presents a recruitment structure based on a respondent-driven sample among 303 heavy drug users from Curitiba collected between July 28, 2008, and October 18, 2009 ([SALGANIK; FAZITO, et al., 2011](#), Web Appendix). Five seeds were chosen within the population, the fifth being a month after the other four since the fourth seed was unsuccessful. Each participant received three coupons and the mean number of recruited individuals per recruiter was around 0.98.

Figure 2 – RDS structure among heavy drug users in Curitiba.



Source: Data extracted from ([SALGANIK; FAZITO, et al., 2011](#)) and figure prepared by the author (2021). The respondents were asked whether they are in any “Centro de Atenção Psicossocial (CAPS)” (Psychosocial Care Center) treatment program for drug use.

### 2.2.2 Assumptions and statistical properties

RDS is a successful recruitment method for reaching hard-to-reach populations since the respondents recruit most of the participants. On the other hand, this characteristic also makes it hard to derive statistical properties without making strong assumptions of the recruitment process. Some hypotheses are related to specific models, which are presented in Section [2.2.3](#):

- sampling is not uniformly random among the individuals since some have more connections than others, which gives them a higher probability of being recruited. Those with more contacts should reduce the weighting in the inferences, but this also relies on another assumption: self-reported degree should be accurately measured ([GILE; HANDCOCK, 2010](#), p. 297);
- recruitment is without replacement, given that respondents are not allowed to participate more than once. It compromises inferences since the probability of inclusion in the survey also depends on the number of individuals participating until the recruitment time ([GILE; HANDCOCK, 2010](#), p. 299). To derive an RDS estimator, [Volz and Heckathorn \(2008](#), p. 81) requires a small sampling fraction to compensate for breaking this assumption;
- homophily* is the tendency of individuals to connect within the same group. For instance, men tend to recruit more men to women. If the

process has zero homophily, it indicates that individuals do not regard the group to recruit. On the other hand, if homophily is one, all the connections are intragroup (HECKATHORN, 2002, p. 20). Heckathorn (2002, p. 21) proved that under certain conditions (see Subsection 2.2.3), the respondent-driven sample is unbiased with respect to homophily if it is equal for each group;

- d) the connections generated by the RDS process item b) violate the independence between the samples through *clustering*, i.e., people are more likely to connect to those similar (AVERY, 2020, p. 14);
- e) respondent-driven sampling produces a branching structure that makes it impossible to observe links between two people who don't recruit each other (GILE; HANDCOCK, 2015, p. 17). It constitutes a missing data problem, according to Crawford (2016, p. 190);
- f) in apparent contraction to item b), to the distribution achieve its convergence and remove the biases induced by the initial sample, enough waves of recruitments are necessary (HECKATHORN, 1997, p. 186);
- g) Goel and Salganik (2009, p. 2225) defines *bottleneck* as the probability of cross-group recruitment. It happens when the recruitment chain remains inside an identified subgroup of individuals. In that situation, “studies should be conducted separately within each tier.” (GILE; BEAUDRY, et al., 2018, p. 75). As an expository example, Toledo et al. (2011, p. S139) observes strong geographical heterogeneity among a population of heavy drug users in Rio de Janeiro.

### 2.2.3 Models for the RDS Process

Since its inception, several authors have tried to better understand and model the RDS process because of its non-probability nature. Each modelling approach aims to approximate even more the network structure to yield more reliable inferences. In this section, we present some of them. Let  $G = (V, E)$  be an undirected graph representing the hidden population, such that  $|V| = N$ , and  $A \in \{0, 1\}^{N \times N}$  its adjacency matrix, where  $A_{ij} = 1$  if there is a connection between individuals  $i$  and  $j$ , and  $A_{ij} = 0$  otherwise. We denote  $|V|$  to mean the number of nodes, and  $|E|$  the number of edges in the graph  $G$ . The choice of an undirected model for the hidden population is very common, but not obliged. The degree of a person is, therefore,  $d_i = \sum_{j=1}^N A_{ij}$ .

Include notation of RDS used posteriorly.

Besides the following models, (GOEL; SALGANIK, 2009) presents RDS as a MCMC.

### 2.2.3.1 First-order Markov process

This approximation was the first model proposed by Heckathorn (1997). He argues RDS recruitment has the characteristic that “any subject’s recruits are a function of his or her type, such as his or her ethnicity; and not of previous events, such as who recruited the recruiter” (HECKATHORN, 1997, p. 182). Consequently, recruitment is modelled as a first-order Markov chain in the space of states generated by the categorical variables, such as ethnicity or gender. The evidence for the above statement is based on chi-square analysis. By these hypotheses, the paper derives three theorems:

**Theorem 2.2.1** (Convergence to equilibrium). *Let  $\{Z_n\}_{n \in \mathbb{N}}$  be the recruitment process. Given that the space space is finite, if the Markov chain is irreducible and aperiodic, then it converges to the stationary distribution and is independent of the initial sample (HECKATHORN, 1997, p. 183).*

*Proof.* A proof is outlined in (LEVIN; PERES, 2017, p. 52-53).  $\square$

**Theorem 2.2.2** (Geometric rate of convergence). *The convergence of the Markov chain generated by RDS recruitment converges to the stationary distribution at a geometric rate (HECKATHORN, 1997, p. 186).*

*Proof.* The same proof given in (LEVIN; PERES, 2017, p. 52-53), demonstrates the geometric convergence.  $\square$

**Theorem 2.2.3** (Unbiased samples). *A respondent-driven sample produces an unbiased sample if all groups have same homophily, that is, the probability of selecting a member within the same group for any group is the same (HECKATHORN, 1997, p. 192).*

*Proof.* Heckathorn (1997, p. 191 - 192) presents a proof for this fact.  $\square$

Heckathorn (2002, p.22) extended this model with the hypothesis that relationships between the individuals are reciprocal. The Random Walk model simplifies this concept by proposing that each recruitment in the social network  $G$  occurs between adjacent nodes with uniform probability and that the process begins with a unique seed. With the assumption that the graph has only one connected component and that the researchers chose the seed with probability proportional to its degree,

Salganik and Heckathorn (2004, p. 209-218) derives sampling probabilities. A proof of asymptotic convergence to the stationary distribution

$$\pi_j^* = \frac{d_j}{\sum_{i=1}^N d_i} \quad (2.6)$$

is provided (SALGANIK; HECKATHORN, 2004, p. 234-235). The authors ponder limitations regarding the validity of these assumptions in real applications and argues that “Empirically checking the reasonableness of the assumptions and further research related to the robustness of the estimation procedure are both problems worthy of further study.” (SALGANIK; HECKATHORN, 2004, p. 230).

### 2.2.3.2 Successive sampling (SS)

The problem with the Random Walk approach with replacement is the assumption of a small sample fraction. It induces biases in prevalence estimates since population size can be small, implying that convergence will not occur or the sample fraction will be high. To adjust for finite population effects, Gile (2011) suggests a successive sampling approach. Along with the sampling, the recruitment probability is proportional to the size of the remaining not recruited population.

The procedure starts sampling an individual  $i$  with probability proportional to degree  $d_i$ . After, it selects another individual with probability proportional without replacement, given by expression (2.7) (GILE, 2011, p. 136).

$$\Pr(G_j = g_j \mid G_1 = g_1, \dots, G_{j-1} = g_{j-1}) = \begin{cases} \frac{d_{g_j}}{2|E| - \sum_{i=1}^{j-1} d_{g_i}}, & g_j \notin \{g_1, \dots, g_{j-1}\} \\ 0, & \text{otherwise,} \end{cases} \quad (2.7)$$

such that  $G_i = g_i$  is the event of the selection of individual  $g_i$  in the step  $i$ . To estimate probabilities, this model assumes that the degree distribution and the population size  $N$  are known (GILE, 2011, Table 2, p. 144), the latter not being necessary to the Random Walk with replacement model.

### 2.2.3.3 Graphical Structure model

Crawford (2016) presented a model to probabilistically reconstruct the subgraph whose nodes are the respondents and edges are their connections. He considered the information brought by the waiting times between recruitments and the remaining coupons with the recruiters to define a probability distribution on the space of subgraphs.

**Definition 2.2.1** (Recruitment graph). The *recruitment graph*  $G_R = (V_R, E_R)$  represents the recruited individuals and the recruitment edges. Therefore  $i \in V_R$  if individual  $i \in V$  was recruited, and  $(i, j) \in E_R$  if individual  $\{i, j\} \in E$  and individual  $i$  recruited individual  $j$ . Notice that  $G_R$  is a *forest*, that is, a collection of trees. (CRAWFORD, 2016, p. 193).

Denote  $n = |V_R|$ . Given that each individual can be sampled only once, it is not possible to observe the *recruitment-induced subgraph*, that is

**Definition 2.2.2** (Recruitment-induced subgraph). It is the induced subgraph  $G_S = (V_S, E_S)$  generated by  $V_R$ , that is,  $V_S = V_R$  and  $\{i, j\} \in E_S$  if  $i, j \in V_R$  and  $\{i, j\} \in E$ . (CRAWFORD, 2016, p. 192).

Denote  $\mathbf{t} = (t_1, \dots, t_n)$  the vector of recruitment times of the individuals such that  $t_1 < \dots < t_n$ , and  $\mathbf{d} = (d_1, \dots, d_n)$  the degrees of the individuals in the same order. Then we define

**Definition 2.2.3** (Coupon matrix). The *coupon matrix*  $C \in \{0, 1\}^{n \times n}$  defined by  $C_{ij} = 1$  if the  $i^{th}$  subject has at least one coupon just before the  $j^{th}$  recruitment event. The row order is the same of  $\mathbf{t}$ . (CRAWFORD, 2016, p. 193).

From the RDS process, the observed data is  $\mathbf{Z} = (G_R, \mathbf{d}, \mathbf{t}, C)$ .

**Definition 2.2.4** (Compatibility). Let  $\hat{G}_S = (\hat{V}_S, \hat{E}_S)$  be an estimate for  $G_S$ . The subgraph  $\hat{G}_S$  is *compatible* with data  $\mathbf{Z}$  if

- a)  $v \in V_R$  if and only if  $v \in \hat{V}_S$ ;
- b)  $\forall (i, j) \in E_R, \{i, j\} \in \hat{E}_S$ ;
- c)  $\forall v \in V_R, \sum_{u \in V_R / \{v\}} \mathbb{1}\{\{u, v\} \in \hat{E}_S\} \leq d_v$ . (CRAWFORD, 2016, p. 197).

We denote  $\mathcal{C}(\mathbf{Z})$  the set of all compatible subgraphs for  $\mathbf{Z}$ .

After the recruitment time  $t_i$ , individual  $i$  is a recruiter until their coupons or non recruited neighbors are exhausted. A node is *susceptible* if it has a link to a recruiter. An edge is susceptible if it connects a recruited and a susceptible node. After  $j$  being recruited, every  $\{i, j\} \in E$  with  $i \in V_R$  is no longer a susceptible edge. Moreover, Crawford (2016, p. 194) assumes that each recruitment time has exponential distribution with parameter  $\lambda$  and it is independent of the recruiter characteristics, neighbors, and all other waiting times. This assumption may fail when homophily is strong. Some interesting propositions follows from this construction (CRAWFORD, 2016, p. 195), but here we focus on  $G_S$ .

Provide an example to explain all the above definitions.

Let  $\tilde{A} \in \{0, 1\}^{n \times n}$  be the adjacency matrix of a compatible estimated subgraph, that is,  $[\tilde{A}]_{ij} = 1$  if and only if  $\{i, j\} \in \hat{G}_S$ . Then

$$[AC]_{ij} = \sum_k [A]_{ik}[C]_{kj} = \sum_k \mathbb{1}(\{i, k\} \in \hat{G}_S \text{ and } k \text{ can recruit in } t_j),$$

that is the number of recruiters connected to  $i$  just before the  $j^{th}$  recruitment, when  $j \leq i$ . Let  $u_i$  be the number of edges linking the sampled node  $i$  with others not sampled. Then,

$$[C^T u]_i = \sum_k [C]_{ki} u_k = \sum_k \mathbb{1}(k \text{ can recruit at } t_i) \cdot \#\text{susceptible edges of } k$$

**Proposition 2.2.1.** *The likelihood of the recruitment times  $w = (0, t_2 - t_1, \dots, t_n - t_{n-1})$  is*

$$L(w|G_S, \lambda) = \left( \prod_{k \text{ isn't seed}} \lambda s_k \right) \exp(-\lambda \mathbf{s}^T w), \quad (2.8)$$

where

$$\mathbf{s} = \text{tril}(\tilde{A}C)^T \mathbf{1} + C^T u$$

indicates the number of susceptible edges just before each recruitment. ([CRAWFORD, 2016](#), p. 197).

*Proof.* A proof of this proposition is given in the online Appendix of ([CRAWFORD, 2016](#)).  $\square$

Setting  $T(\tilde{A}) = -\lambda \mathbf{s}$  and  $B(\tilde{A}) = \sum_{k \text{ isn't seed}} \log(\lambda s_k)$ , the likelihood from above can be normalized to obtain the probability

$$P(\tilde{A}|w) \propto \exp \left[ T(\tilde{A})^T w + B(\tilde{A}) \right]$$

which can be interpreted as an Exponential Random Graph Model (ERGM) ([CRAWFORD, 2016](#), p. 198). Finally, from a Bayesian perspective (see Section 2.5), one can define prior distributions over  $G_S$  and  $\lambda$  to obtain,

$$p(G_S, \lambda|G_R, C, d, t) \propto L(w|G_S, \lambda) \pi(G_S, \lambda), \quad (2.9)$$

where  $\pi(G_S, \lambda)$  is a prior density. A Metropolis-within-Gibbs sampling scheme is used to draw pairs  $(G_S, \lambda)$ . A simulated annealing procedure can be used to obtain a sequence that converges to the maximum a posteriori. An application of this model was the estimation of the hidden population size with the additional assumption that the graph  $G$  has Erdős-Rényi distribution ([CRAWFORD; WU; HEIMER, 2018](#)).

### 2.2.3.4 New model

Maybe only a superficial explanation about the model, since it is really difficult, despite being apparently very nice. ([MCLAUGHLIN, 2021](#))

## 2.2.4 Prevalence estimators

In this subsection, we outline five very common RDS proportion estimators presented in the literature based on the modelling from Subsection 2.2.3. They are apparent prevalence estimators and can be used for prevalence estimate through equation (2.5) in a frequentist approach. Then:

- a) *naive estimator*: it is the sample proportion

$$\hat{\theta}_{\text{naive}} = \frac{1}{n} \sum_{i=1}^n y_i,$$

as in equation (2.4);

- b) *Salganik-Heckathorn (SH) RDS estimator*: Considering the Random Walk approximation, [Salganik and Heckathorn \(2004\)](#) built this estimation regarding the sampling probabilities. Let  $N_T = \sum_{i \neq j} A_{ij}y_i(1 - y_j)$  be the number of connections between individuals with and without the disease,  $\bar{d}_1 = \frac{\sum_{i=1}^N \sum_{j \neq i} A_{ij}y_i}{\sum_{i=1}^N y_i}$  the mean degree of ill individuals,  $\bar{d}_0$  the mean degree of not ill individuals with a similar formula, and  $N_1 = N\theta$ . [Salganik and Heckathorn \(2004, p. 218\)](#) derives that

$$\theta = \frac{\bar{d}_0 c_{01}}{\bar{d}_0 c_{01} + \bar{d}_1 c_{10}},$$

where

$$c_{01} = \frac{N_T}{(N - N_1)\bar{d}_0} \text{ and } c_{10} = \frac{N_T}{N_1 \bar{d}_1},$$

and that

$$\hat{\theta}_{\text{SH}} = \frac{\hat{d}_0 \hat{c}_{01}}{\hat{d}_0 \hat{c}_{01} + \hat{d}_1 \hat{c}_{10}}, \quad (2.10)$$

is the prevalence estimator, such that  $\hat{d}_0$ ,  $\hat{d}_1$ ,  $\hat{c}_{01}$ , and  $\hat{c}_{10}$  are estimated for these quantities.

- c) *Volz-Heckathorn RDS (VH) estimator*: With similar assumptions to the previous one, [Volz and Heckathorn \(2008, p. 85\)](#) shows that the inclusion probability of individual  $i$  in the sample is  $\pi_i \propto d_i$  and the corresponding proportion estimator is

$$\hat{\theta}_{\text{VH}} = \frac{\sum_{i=1}^n y_i d_i^{-1}}{\sum_{i=1}^n d_i^{-1}}. \quad (2.11)$$

The assumptions for  $\hat{\theta}_{\text{VH}}$  were highlighted in Subsubsection 2.2.3.1, and are summarized in (Table 1 GILE; BEAUDRY, et al., 2018, p. 71).

- d) *Successive sampling (SS) estimator:* Under the successive sampling approximation for RDS, Gile (2011, p. 137-138) derives an estimate considering the without replacement assumption. It is of the form

$$\hat{\theta}_{\text{SS}} = \frac{\sum_{i=1}^n y_i w_i}{\sum_{i=1}^n w_i}, \quad (2.12)$$

where  $w_i$  is calculated algorithmically, taking account the finite population effect. If the sampling fraction is small, this estimator is similar to VH estimator. Otherwise, when it grows, VH is biased. The limitation of SS estimator is that  $N$  is assumed to be known, which is rarely the case. Gile (2011, p. 140) did a sensitivity analysis on population size estimate.

- e) *RDS-B estimator:* (BASTOS; BASTOS, et al., 2018) proposes a pseudo-posterior approach to estimate prevalence. Let

$$Y_i \sim \text{Bernoulli}(\theta_i) \text{ with } \text{logit}(\theta_i) = \alpha,$$

where logit is explained in Section 2.3. Defines  $\delta_i \propto n \cdot d_i^{-1}$  such that  $\sum_{i=1}^n \delta_i = n$ , based on the weights suggested by Volz and Heckathorn (2008). The pseudo-likelihood is written as follows:

$$L(\alpha \mid Y = y) = \prod_{i=1}^n \Pr(Y_i = y_i \mid \alpha)^{\delta_i}.$$

In a Bayesian perspective (see Section 2.5), inferences are based on the posterior distribution and Bastos, Bastos, et al. (2018, p. S18) used weakly informative priors for  $\alpha$ . In this case, a pseudo-posteriori is used. This estimator has the advantage of allowing prior information as convenient, but it suffers from the same limitations as VH and SH estimators, since the weights are derived from a Random Walk approximation.

Ott et al. (2019) and Fellows (2019) extended these estimators. The former presented a similar estimator to SH estimator, yet more robust. The latter introduced homophily into the model. Besides these estimators, Avery (2020) suggested binary logistic regression methods and other extensions through Generalized Linear Models (see Section 2.3).

## 2.2.5 Regression methods

According to Gile, Beaudry, et al. (2018, p. 86), “RDS suffers from two particular challenges for multivariate modeling: unknown sampling weights and

unknown dependence structure.” These two problems led to different approaches in the literature. Avery (2020, p. 13-15) has a good review on the topic. Spiller (2009) suggests to model dependence as mixed effects. Bastos, Pinho, et al. (2012) performs a binary regression to prevalence estimation through a hierarchical model where correlation structure was modelled as a Conditionally autoregressive (CAR) model (see Section XXX). Yauck et al. (2021) includes homophily in a similar model, but with a Simultaneous Autoregressive (SAR) model (CITATION) for correlation.

It would be nice to cite Camila, but reference not found.

### 2.2.6 Bootstrap methods for uncertainty quantification

(BARAFF; MCCORMICK; RAFTERY, 2016), (SALGANIK, 2006).

### 2.2.7 Diagnosis of RDS

(GILE; JOHNSTON; SALGANIK, 2015)

## 2.3 Generalized linear models

Let  $\mathbf{y} \in \mathbb{R}^n$  be a realization of a random variable  $Y : \Omega \rightarrow \mathbb{R}^n$  associated with a phenomena such that each component  $Y_i$  is independent of the others. Set  $\mu = \mathbb{E}[Y]$ . The classical linear model assumes that  $Y_i \stackrel{iid}{\sim} \text{Normal}(\mu_i, \sigma^2)$  and  $\mu = \mathbf{X}\beta$ , such that  $\beta \in \mathbb{R}^k$  is an unknown parameter vector and  $\mathbf{X} \in \mathbb{R}^{n \times k}$  is the data, where  $\mathbf{X}_{ij}$  is the measure of the  $j$ -th covariate in the  $i$ -th individual. Non constant variance for each  $Y_i$  is a possible variation for this model.

Generalized linear models (GLM) extend the above model. In order to understand this extension, we follow McCullagh and Nelder (2019, p. 27) setting

$$\eta = \mathbf{X}\beta \quad \text{and} \quad \eta_i = g(\mu_i), i = 1, \dots, n,$$

such that  $g(\cdot)$  is a monotonic differentiable function and is named *link function*. Therefore, “the link function relates the linear predictor  $\eta$  to the expected value  $\mu$  (McCULLAGH; NELDER, 2019, p. 31). Notice that in the classical linear model,  $g$  is the identity function, but it can be generalized. Another possible generalization is the distribution of  $Y$ , which may be any from the Exponential Family distribution (ROBERT, 2007, p. 115).

When  $Y_i$  has Bernoulli distribution with probability of success  $\mu \in (0, 1)$ , the link function must have its image over the open interval  $(0, 1)$  and domain in the real line. The classical are the following:

- a) *logit*:  $\eta = \log(\mu/(1 - \mu))$  that represents the log odds of  $Y_i = 1$ ;
- b) *probit*:  $\eta = \Phi^{-1}(\mu)$  where the  $\Phi(\cdot)$  is the Normal cumulative distribution function;
- c) *complementary log-log*:  $\eta = \log(-\log(1 - \mu))$ .

This work focus on Logistic regression, which is the most common inferencial procedure for binary response, such as having or not a disease.

## 2.4 Conditionally autoregressive models

This construction follows the Gaussian case ([BANERJEE; CARLIN; GELFAND, 2003](#), Section 3.3.1).

## 2.5 Bayesian statistics

We can represent our beliefs and information about unknown quantities through probabilities. There are two more common interpretations: frequentist and Bayesian. While the frequentists define probability as the limit of a frequency in a large number of trials, the Bayesians represent an individual's degree of belief in a statement that is updated given new information. This philosophy allows assigning probabilities to any event, even if a random process is not defined ([STATISTICAT, 2016](#)).

In 1761, Reverend Thomas Bayes wrote for the first time the Bayes' formula relating the probability of a parameter after observing the data with the evidence (written through a likelihood function) and previous information about the parameter. Pierre Simon Laplace rediscovered this formula in 1773 ([ROBERT, 2007](#)), and this theory became more common in the 19th century. After some criticisms, a modern treatment considering Kolmogorov's axiomatization of the theory of probabilities started after Jeffreys in 1939. The recent development of new computational tools brought these ideas again.

Therefore, Bayesian inference is the process of inductive learning using Bayes' rule, where inductive means that characteristics of a population are learned from a subset of it. We generally express numerical characteristics of the population as a parameter  $\theta$  which is indirectly observed through numerical descriptions  $y$  of the population. Both are uncertain until the observation of a sample, when its information can decrease our uncertainty about the population characteristics ([HOFF, 2009](#), p. 1-2).

The set of all possible outcomes  $y$  forms the *sample space*  $\mathcal{Y}$ , while the set of all possible parameters forms the *parameter space*  $\Theta$ . Bayesian inference is composed by the following:

- a) *prior distribution*: A probability distribution defined over  $\Theta$  that quantifies our beliefs about  $\theta$  before observing the data;
- b) *sampling model*: A probability distribution of the data generation process that express our belief that  $y \in \mathcal{Y}$  is the outcome when  $\theta \in \Theta$  is true. When it is seen as function of the parameter, it is called *likelihood function*;
- c) *loss function*: Only in a decision theory framework, it measures the error of a estimative  $\delta \in \Theta$  in comparison to  $\theta$ ;
- d) *posterior distribution*: Once we get the data  $y$ , it represents our updated beliefs out the parameter conditioned All inferences are based on this probability distribution.

Bayes' theorem establishes that when the sampling model is absolutely continuous with respect to some measure  $\nu$  with conditional density  $f_{Y|\theta}(y | \theta)$  and the prior distribution is a well defined probability measure  $\mu_\theta$ , the posterior distribution  $\mu_{\theta|Y}(\cdot | y)$  is absolutely continuous with respect to  $\mu_\theta$  almost surely and its Radon-Nikodym derivative is (SCHERVISH, 2012, p. 16)

$$\frac{d\mu_{\theta|Y}}{d\mu_\theta}(\theta|y) = \frac{f_{Y|\theta}(y | \theta)}{\int_{\Theta} f_{Y|\theta}(y | t) d\mu_\theta(t)}. \quad (2.13)$$

When the prior distribution is absolutely continuous with respect to the Lebesgue measure, equation (2.13) resumes to

$$p(y|\theta) = \frac{f(y | \theta)\pi(\theta)}{\int_{\Theta} f(y | t)\pi(t) dt}. \quad (2.14)$$

## 2.6 Computational methods

### 2.6.1 Hamiltonian Monte Carlo

We follow (BETANCOURT, 2017). This method was developed in the late 1980s as Hybrid Monte Carlo to tackle calculations in Lattice Quantum Chromodynamics. Instead of moving in the parameter space randomly with uninformed jumps, the direction from the vector field given by the gradients are used to trace out a trajectory through the \*typical set\*, the region which has significant contribution to the expectations. However, if only the gradient was used, the trajectory would pull towards the mode of the distribution, so more geometric constraints are needed. In

order to a satellite rotate around the Earth, we have to endow it with enough momentum to counteract the gravitational field, turning the system into a conservative one.

First, we introduce auxiliary momentum parameters  $p_n$  (lift) of the same dimension from the parameter space  $\Omega \subseteq \mathbb{R}^D$ . Then  $q_n$  turns to  $(q_n, p_n)$ , with the use the joint probability distribution  $\pi(q, p) = \pi(p | q)\pi(q)$ . Particularly, we use

$$\pi(q, p) = e^{-H(q, p)},$$

such that  $H$  is the \*Hamiltonian\*. Note that  $H(q, p) = -\log \pi(p | q) - \log \pi(q) =: K(p, q) + V(q)$ . We call  $K$  the kinetic energy, and  $V$  the potential energy. The vector field is generated by Hamilton's equations,

$$\begin{aligned}\frac{dq}{dt} &= \frac{\partial H}{\partial p} = \frac{\partial K}{\partial p} \\ \frac{dp}{dt} &= -\frac{\partial H}{\partial q} = -\frac{\partial K}{\partial q} - \frac{dV}{dq}.\end{aligned}$$

Therefore, we are able to define the Hamiltonian flows  $\phi_t : (p, q) \rightarrow (p, q), \forall t \in \mathbb{R}$ .

### 2.6.1.1 Diagnostics

The importance of diagnosing. The potential problems that it can show.

- Divergent transitions;
- Transitions that hit the maximum tree depth;
- Low E-BFMI values;
- Low effective samples sizes;
- $\hat{R} \notin (0.95, 1.05)$ .

## 2.6.2 Metropolis-within-Gibbs

If this method is used.

### 3 Prevalence modelling and regression methods

Fisher (1922, p. 311) stated that the objective of statistics is to reduce the data since its volume is impossible to comprehend by the researchers. In that sense, few parameters should represent the whole phenomenon catching the most relevant information. Years later, Newman studied the theory of modelling which can be divided in three aspects (LEHMANN, 2012, p. 161):

- a) models of complex phenomena are created by putting together simple building elements that the researcher is familiar with and can handle;
- b) there are two types of models: the *explanatory models*, which will be focused on this work, and the *interpolatory formulae*.
- c) An explanatory theory necessitates a thorough understanding of the scientific context of the problem. In this regard, we investigated questions involving Respondent-driven sampling and prevalence estimation as introduced in Chapter 2.

In this chapter, we develop models that enclose these ideas building each block separately. For a Bayesian modelling, we assume that each parameter of the model has a probability distribution that incorporates the researcher's uncertainty about it. For each individual, we observe  $k$  covariates that are possible risk factors represented by the vector  $\mathbf{x}_i \in \mathbb{R}^k$  of the  $i^{th}$  individual. We denote  $\theta_i$  the probability of the  $i$ -th individual have been exposed to the disease that depends on the prevalence  $\theta$  and  $\mathbf{x}_i$ . We also consider the dependence of sampling from RDS as a spatial random effect. The probability of positive test in the  $i^{th}$  individual is denoted by  $p_i$ .

Another important feature of the model is that sensitivity and specificity have the same distribution for all individuals and it only depends on the test used to diagnose. This is an assumption that must be analysed for each particular case. For instance, COVID-19 Sofia test has different sensibility and specificity for symptomatic and asymptomatic individuals (Table 1 MITCHELL et al., 2021, p. 3).

From above, we develop three different models: the first considers perfect tests, that is,  $\gamma_s = \gamma_e = 1$  and no spatial random effect; the second considers imperfect tests, regarding  $\gamma_s$  and  $\gamma_e$ , but ignoring the RDS structure; and the third one has imperfect tests and RDS structure. Some considerations are made to improve the

model's limitations.

The implementation of the following models were in the statistical computation platform Stan ([CARPENTER et al., 2017](#)) within Python Interface PyStan ([RIDDELL; HARTIKAINEN; CARTER, 2021](#)) which uses an implementation for HMC algorithm. All the codes are written in [Appendix B](#).

## 3.1 Perfect tests

The first model supposes the samples are independent and the test is perfect, which means that  $\theta_i = p_i$  for all  $i$ . Therefore it only considers the risk factors  $\mathbf{x}_i$ .

$$\begin{aligned} Y_i \mid \theta_i &\sim \text{Bernoulli}(\theta_i), \\ g(\theta_i) &= g(\theta) + \mathbf{x}_i^T \beta, \end{aligned} \tag{3.1}$$

where  $g(\cdot)$  is the logit function. The parameter  $\beta \in \mathbb{R}^k$  is the risk effects. For Bayesian inference, priors on  $\beta$  and  $\theta$  must be included. We use  $\beta \sim \text{Normal}(\mu_\beta, \Sigma_\beta)$  and  $\theta \sim \text{Beta}(a^p, b^p)$ , where the vector  $\mu_\beta \in \mathbb{R}^k$ , the symmetric positive-definite matrix  $\Sigma_\beta \in \mathbb{R}^{k \times k}$ , and the positive real values  $a^p, b^p \in \mathbb{R}_{>0}$  are fixed hyperparameters. Inferences about  $\beta$  and  $\theta$  are based on the posterior distribution. Keeping the notation of Section [2.3](#), we denote  $\mathbf{X}$  the covariate matrix.

*Remark 3.1.1* (Interpretation of prevalence). According to the model formulation, if the risk factors are zero, i.e  $\mathbf{x}_i = 0$ , the probability of the  $i$ -th individual having been exposed is the prevalence  $\theta$ , which means that in a population with no risk effects, the probability of a person having the disease is exactly the proportion in this population.

### 3.1.1 Identifiability

A formal definition for identifiability regards the likelihood function ([XIE; CARLIN, 2006, p. 3459](#)):

**Definition 3.1.1.** Let  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$  be the family of probability distributions for  $\mathcal{Y}$ . This model is *identifiable* if for any  $\theta', \theta'' \in \Theta$ ,

$$\forall y \in \mathcal{Y}, P_{\theta'}(Y = y) = P_{\theta''}(Y = y) \implies \theta' = \theta''.$$

The family distribution from model (3.1) is the logistic regression parametrized by  $(\theta, \beta)$  and conditioned on observing the regressor  $\mathbf{X}$ , with  $\mathcal{Y} = \{0, 1\}^n$ . Defining

$\beta_0 = g(\theta)$ , we may rewrite it as

$$Y_i \mid \tilde{\beta}, \tilde{\mathbf{x}}_i \sim \text{Bernoulli}(g^{-1}(\tilde{\mathbf{x}}_i^T \tilde{\beta})),$$

such that  $\tilde{\beta}$  concatenate  $\beta_0$  and  $\beta$ , and  $\tilde{\mathbf{x}}_i$  concatenate 1 and  $\mathbf{x}_i$ . Küchenhoff (1995, p. 7) gives a formal proof for the identifiability of this representation.

In the Bayesian paradigm, inferences are based on the posterior distribution. Therefore, identifiability should consider the prior distribution. Lindley (1972, p. 46) argued that proper priors are sufficient to handle identifiability problems in the Bayesian perspective, which means that a well-defined posterior probability distribution is enough for parameter identification. A formal definition for *Bayesian identifiability* is the following: if  $p(\theta \mid \beta, y, \mathbf{X}) = p(\theta \mid \beta)$ , the data  $y$  is uninformative for  $\theta$  when  $\beta$  is known. The definition is analogous if  $\beta$  and  $\theta$  change places. However, Gelfand and Sahu (1999, p. 248) proved that this definition is equivalent to likelihood identifiability.

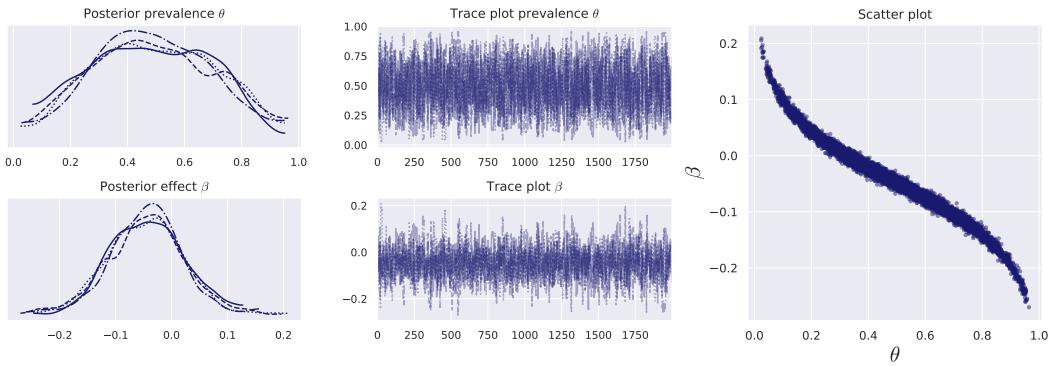
Despite the identifiability of the model, it may be hard to sample from the posterior distribution depending on the value of  $\mathbf{x}$ . As an example, consider the following experiment:

- (i) generate 500 covariates  $X_i \sim \text{Normal}(15, 1)$ ;
- (ii) let  $\beta = 0.1$ ,  $\theta = 0.1$ , and  $\theta_i = g^{-1}(g(\theta) + X_i\beta)$  for  $1 \leq i \leq 500$ ;
- (iii) for each  $i$ , sample  $Y_i \sim \text{Bernoulli}(\theta_i)$ ;
- (iv) let  $a^p = 1$ ,  $b^p = 1$ ,  $\mu_\beta = 0$ , and  $\Sigma_\beta = 1$  the hyperparameters for the prior distributions (weakly informative);
- (v) make 1000 warm-up and 1000 sampling iterations using Stan given the data  $(Y_1, X_1), \dots, (Y_n, X_n)$ .
- (vi) make 2000 warm-up and 2000 sampling iterations using Stan given the data  $(Y_1, X_1), \dots, (Y_n, X_n)$ .

The HMC sampler took around 8.39s. Figure 3 presents the results through the posterior distribution, the trace plot, and the strong posterior correlation between  $\theta$  and  $\beta$ . To address this problem, subtracting the mean  $\bar{x}$  is a default procedure (OGLE; BARBER, 2020, p. 5). After centering the data around the mean, the HMC sampler took around 1.39s, and the improved results are shown in Figure 4.

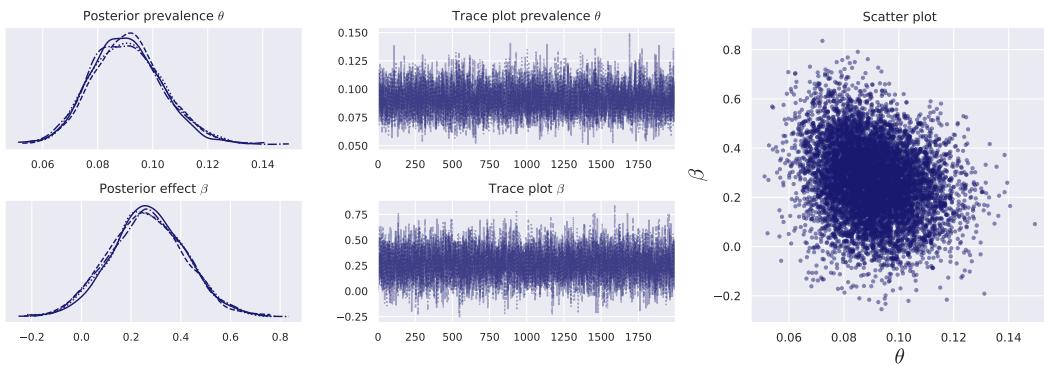
We observe that the interpretation of prevalence from Remark 3.1.1 changes from centred and uncentered since the meaning of  $\mathbf{x}_i = 0$  is different. Along with

Figure 3 – Posterior distribution, trace plot, and posterior samples of parameters  $\theta$  and  $\beta$  from model (3.1) with uncentered covariate.



Source: Prepared by the author (2021).

Figure 4 – Posterior distribution, trace plot, and posterior samples of parameters  $\theta$  and  $\beta$  from model (3.1) with centralized covariate.



Source: Prepared by the author (2021).

In this discussion, it is usual to divide the centred variable by its standard deviation, to put all predictors on a common scale. Discussions about the problems caused by standardizing are outside the scope of this work. Gelman (2008) suggests to divide continuous variables by 2 times the standard deviation to allow “the coefficients to be interpreted in the same way as binary deviation.” (GELMAN, 2008, p. 2867) Binary inputs are not standardized since their coefficients are easily interpretable.

Other identifiability problems arising from the input variables are collinearity and *separation* (GELMAN; JAKULIN, et al., 2008, p. 1360-1361). The latter occurs if a linear combination of a subset of the predictors gives a perfect prediction for the binary outcome. For instance, when a linear combination of the predictors is greater than a threshold if and only if  $y = 1$ .

### 3.1.2 Simulated data

To present a sanity check about the functionality of model (3.1) and to validate the properties of the estimation procedure, we simulate fake data from the model and make inferences about the result. We follow the experiment from Section 3.1.1. Table 1 summarizes the experiment parameters.

Table 1 – Experiment settings for the simulation of model (3.1).

Exp	$n$	$k_c$ (normal)	$k_c$ (cauchy)	$k_b$	$\beta$	$\theta$
1	100	3	0	2	[-0.1, 2.5, 1.4, -1.8, 0.3]	0.05
2	100	3	0	2	[-0.1, 2.5, 1.4, -1.8, 0.3]	0.9
3	100	2	2	1	[-0.1, 2.5, 1.4, -1.8, 0.3]	0.1
4	5000	40	5	5	$F$ distribution	0.1

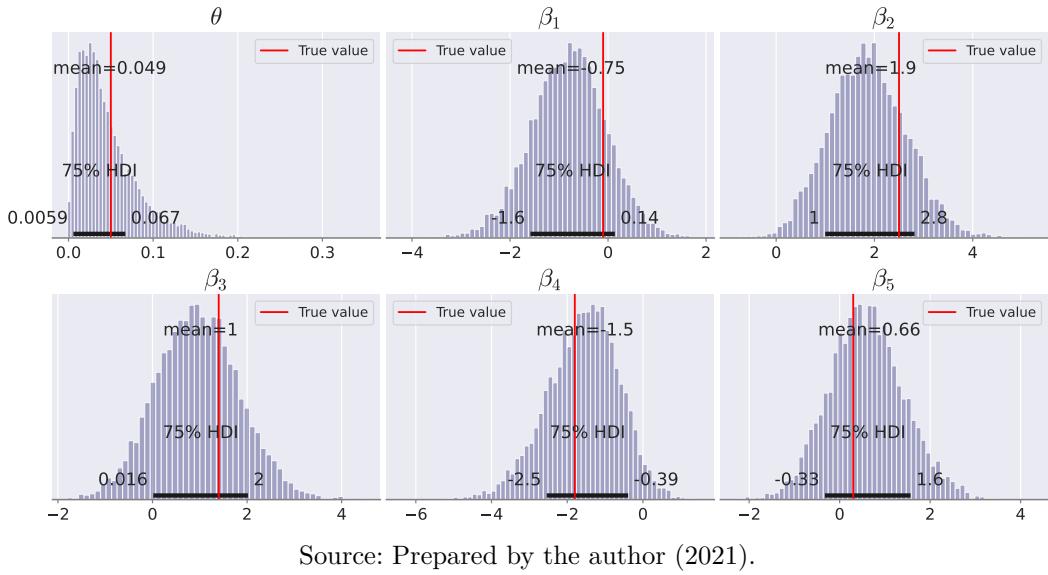
Source: Prepared by the author (2021). We denote  $n$  for number of samples,  $k_c$  for the number of continuous variables, and  $k_b$  for binary variables. Between paranthesis, *normal* means that the variables were generated from a Multivariate Normal with prespecified parameters, and *cauchy* from a Cauchy distribution.  $F$  distribution is  $\text{Normal}(\mu = 0, \sigma = 2)$  with probability 0.3, and 0 otherwise.

We primally look at the settings from experiment 1. With a non-informative prior for  $\theta$  (Jeffreys prior  $\text{Beta}(1/2, 1/2)$ ) and a weakly informative for  $\beta$  (zero mean and covariance matrix four times the identity matrix), Figure 5 shows the posterior distributions for the parameters. The prevalence estimate is good despite Jeffreys' prior. When the distance between the prior and the true value is large, the inferences seem to be biased. However, this makes sense regarding the model. For instance, for  $\beta_2$ , before observing the data, we put 0.7 mass probability for values lesser than 0.1. The data decreased it to 0.125. This highlights the importance of a well defined prior distribution. The values for Bulk ESS was greater than 3000 for all parameters, while Tails ESS were greater than 2200 with 1000 warmup and 1000 sampling iterations, and 4 chains. For all parameters  $\hat{R} = 1$ . Trace plots and scatter plots were also good and we omit here since they do not bring new information for the discussion.

Although we are performing Bayesian inference, frequentist properties can be accessed through simulation. After 1000 simulations varying the input data  $Y$ , the 75% credible interval included the true parameters in 75.8%, 78.8%, 76.4%, 77.5%, 67.3%, and 72.2% of the times, respectively for  $\theta, \beta_1, \dots, \beta_5$ . Each simulation had 100 samples and weakly informative priors for  $\beta$  and  $\theta$ . Figure 6 compares the predicted and simulated probabilities. The green area is delimited by the curves generated by  $2\sqrt{\theta_i(1 - \theta_i)/n}$ , where  $n = 500$  is the number of points. It was increased to show a larger variety of points. This area is a  $\pm 2$  standard-error bounds.

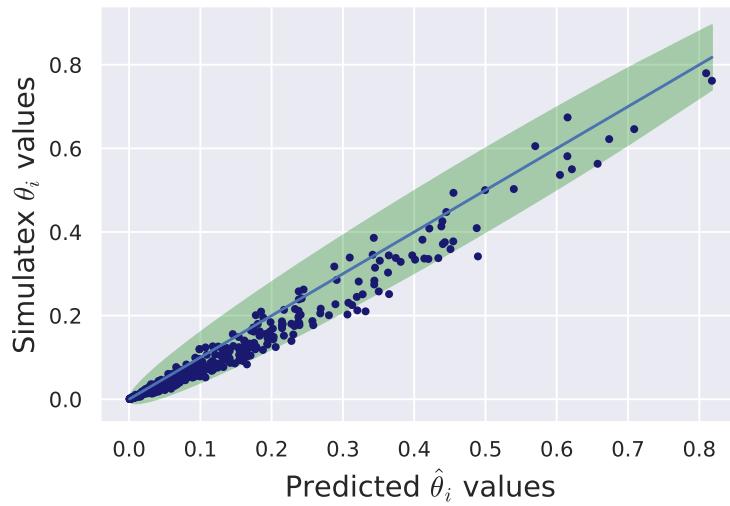
Experiment 2 is used to see if these properties repeat when the prevalence

Figure 5 – Posterior distribution for parameters of model (3.1).



Source: Prepared by the author (2021).

Figure 6 – Comparing predicted and simulated probabilities of having the disease from model (3.1).

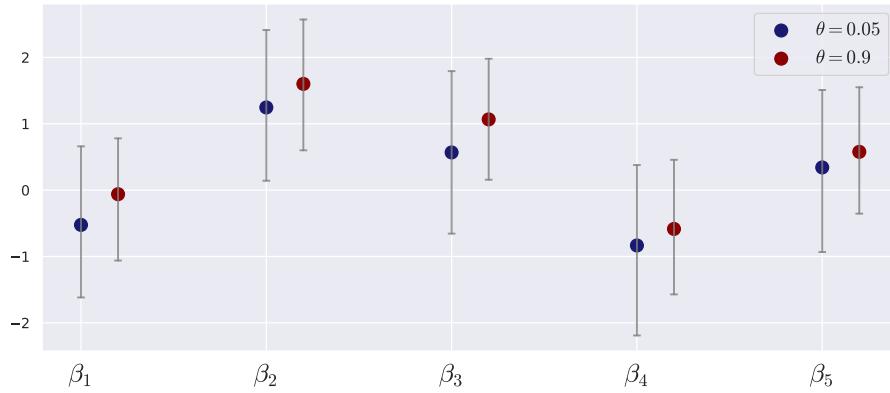


Source: Prepared by the author (2021).

is higher. The same regressors were used for the comparison, but the input data  $Y$  were generated with different prevalences. With prevalence being 0.9, the estimates were a little high for all coefficients as Figure 7 presents. This is related to the fact that the posterior mean underestimated the true value for this experiment. After increasing the number of samples, the estimates were closer, as expected.

The third experiment aims to analyse what happens if some covariates have a heavier tail. No big difference was noticed despite the existence of some

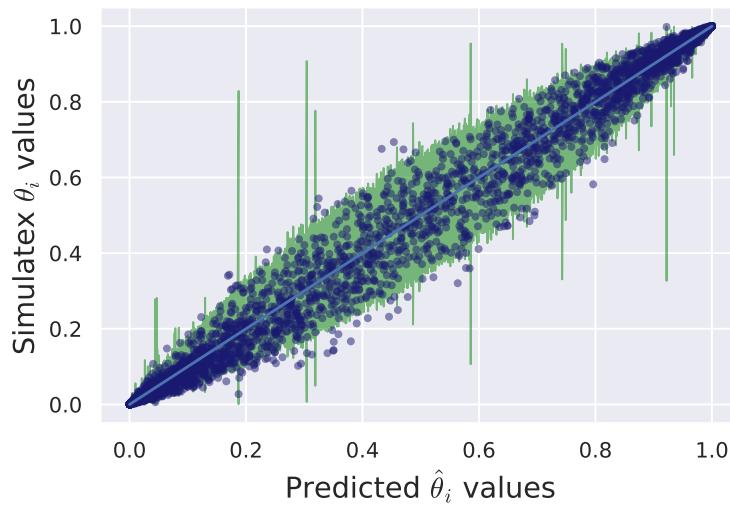
Figure 7 – Comparing posterior mean and 94% credibility intervals for  $\beta$  in model (3.1) with the same regressors  $\mathbf{X}$  but different prevalences.



Source: Prepared by the author (2021).

individuals very different from the others. At last, the fourth experiment increases the dimensionality to observe the number of effective samples. Each chain took around 3 minutes, instead of the 3s needed for the previous experiments. From the 51 parameters, 48 had the true values in the 95% HDI credible interval. The Bulk ESS was greater than 4500 for 95% of the parameters. Figure 8 presents how the predicted probabilities for each individual behaves in this case. Here the green area is delimited by the 95% credible interval.

Figure 8 – Comparing predicted and simulated probabilities of having the disease from model (3.1) with high dimension.



Source: Prepared by the author (2021).

## 3.2 Sensitivity and specificity

In this section, we describe a model for estimating the sensitivity and specificity of a diagnostic test. This model is relevant to analyze and experiment with different prior specification approaches. Suppose having a gold standard test and another test, for instance, a simpler, faster, or less invasive one, which we want to estimate the accuracy by the sensitivity and specificity. In this scenario, true positive (negative) individuals are those who tested positive (negative) by the gold standard. Therefore, in a population with  $n_{\gamma_s}$  true positives and  $n_{\gamma_e}$  true negatives, we denote

$$\begin{aligned} y_{\text{pos}} \mid \gamma_s &\sim \text{Binomial}(n_{\gamma_s}, \gamma_s), \\ y_{\text{neg}} \mid \gamma_e &\sim \text{Binomial}(n_{\gamma_e}, \gamma_e), \end{aligned}$$

such that  $y_{\text{neg}}$  are negative tests on known negative subjects and  $y_{\text{pos}}$  are positive tests on known positive. In the Two-by-two formulation from [Chart 1](#), we have

Chart 2 – Two-by-two table with the model specification.

	$Y = 0$	$Y = 1$	Total
$Y^{\text{true}} = 0$	$y_{\text{neg}}$	$n_{\gamma_e} - y_{\text{neg}}$	$n_{\gamma_e}$
$Y^{\text{true}} = 1$	$n_{\gamma_s} - y_{\text{pos}}$	$y_{\text{pos}}$	$n_{\gamma_s}$
Total	$n_{\gamma_s} + y_{\text{neg}} - y_{\text{pos}}$	$n_{\gamma_e} + y_{\text{pos}} - y_{\text{neg}}$	$n_{\gamma_s} + n_{\gamma_e}$

Source: Prepared by author (2021).

In Bayesian analysis, we have to define a prior distribution with density  $\pi$  for the parameters  $(\gamma_e, \gamma_s)$ . For this, we consider three different approaches:

- a) prior distributions are specified independently for each parameter and each one has a beta distribution, i.e,

$$\pi(\gamma_e, \gamma_s) = \pi(\gamma_e)\pi(\gamma_s) \propto \gamma_s^{a_s}(1 - \gamma_s)^{b_s}\gamma_e^{a_e}(1 - \gamma_e)^{b_e},$$

for  $a_s, b_s, a_e$ , and  $b_e$  being pre-determined positive real hyperparameters;

- b) bivariate normal distribution in the log odds space, i.e,

$$(\text{logit}(\gamma_e), \text{logit}(\gamma_s)) \sim \text{Normal}(\mu_\gamma, \Sigma_\gamma),$$

such that the vector  $\mu_\gamma \in \mathbb{R}^2$  and the covariance matrix  $\Sigma_\gamma \in \mathbb{R}^{2 \times 2}$  are pre-determined hyperparameters;

- c) a bivariate beta distribution described in [Appendix A](#) with parameters  $\alpha_1, \dots, \alpha_4 \in \mathbb{R}_{>0}$ .

### 3.2.1 Independent beta distribution priors

If the knowledge of the specificity affects the range of most possible values of the sensitivity, or vice-versa, there is antecedent information about the correlation between the parameters. When this is not the case, a possible independent prior formulation is the usage of Beta distribution since it is bounded in the interval  $[0, 1]$  and it is reasonably flexible in its shape. Another good reason for this choice is that the beta distribution forms a conjugate family with the binomial distribution, that is, if the likelihood has binomial distribution and the prior has beta distribution, then the posterior has beta distribution.

When considering separated experiments for specificity and sensitivity, there is no information about their correlation, which is the case for our model. Then we define the the prior distributions

$$\begin{aligned}\gamma_e &\sim \text{Beta}(a_e, b_e), \\ \gamma_s &\sim \text{Beta}(a_s, b_s), \\ \theta &\sim \text{Beta}(a_\theta, b_\theta).\end{aligned}$$

Using data from (BENNETT; STEYVERS, 2020) about COVIDPrior information of these quantities lead to a bivariate analysis cite:t‘guo2017bayesian’. As we have already mentioned, the definitions of \*sensitivity\* and \*specificity\* can be expressed as below: -19 seroprevalence in Santa Clara:

$$\begin{aligned}y/n &= 50/3330, \\ y_{\text{neg}}/n_{\gamma_e} &= 399/401, \\ y_{\text{pos}}/n_{\gamma_s} &= 103/122,\end{aligned}$$

we fit the model and obtain the results showed in Figure 9. All the codes were done in *Stan* and *PyStan*.

### 3.2.2 Hierarchical partial pooling prior

Other approach considers more than one study about specificity and sensitivity. A *hierarchical partial pooling* model for these studies can be done in the following way:

$$\begin{aligned}\text{logit}(\gamma_s^j) &\sim \text{Normal}(\mu_{\gamma_s}, \sigma_{\gamma_s}), \\ \text{logit}(\gamma_e^j) &\sim \text{Normal}(\mu_{\gamma_e}, \sigma_{\gamma_e}),\end{aligned}$$

for  $1 \leq j \leq K$  studies, such that the first study is the considered one. Partial pooling because the parameters can be sampled from the same distribution. Hierarchical

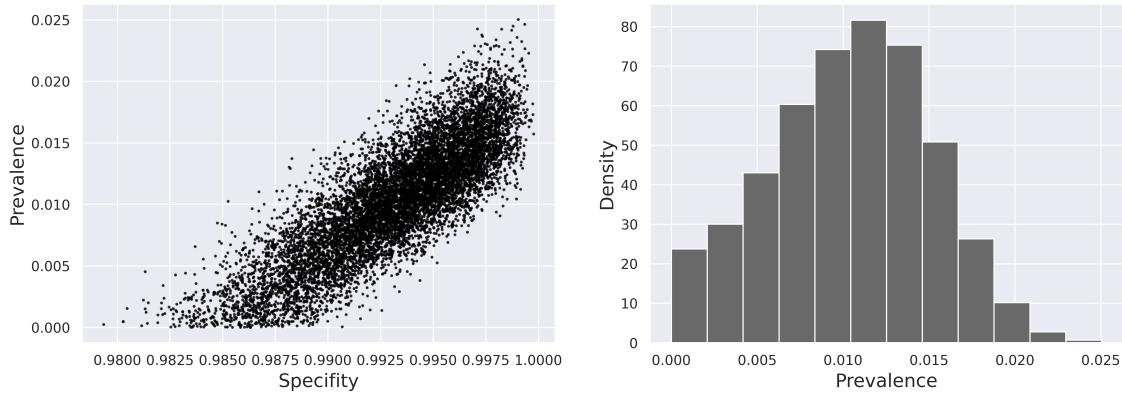


Figure 9 – Scatter plot of posterior simulations of prevalence against specificity and histogram of posterior simulations of the prevalence.

because the parameters of this distribution have its one prior distributions. For instance,

$$\begin{aligned}\mu_{\gamma_s} &\sim N(0, 10), \\ \mu_{\gamma_e} &\sim N(0, 10), \\ \sigma_{\gamma_s} &\sim N^+(0, 1), \text{ and} \\ \sigma_{\gamma_e} &\sim N^+(0, 1),\end{aligned}$$

where  $N^+(a, b)$  is the truncated normal distribution in  $[0, +\infty)$ .

### 3.2.3 Bivariate Beta prior

Finally, we studied a joint distribution for specificity and sensitivity, a possible bivariate beta distribution built in ([OLKIN; TRIKALINOS, 2015](#)). This distribution is derived from a Dirichlet distribution of order four. Let  $U = (U[1], \dots, U[4]) \sim \text{Dirichlet}(\boldsymbol{\alpha})$ , where  $\boldsymbol{\alpha} \in \mathbb{R}_+^4$ . Therefore, defining  $X = U[1] + U[2]$  and  $Y = U[1] + U[3]$ , we will have that  $(X, Y)$  has a well-defined probability distribution in  $[0, 1] \times [0, 1]$  such that  $X$  and  $Y$  have marginally beta distributions, and they have correlation in all space. Depending on the definition of  $\boldsymbol{\alpha}$ , the correlation between the variables range from -1 and 1. Figure 10 shows some examples of this construction.

In this section, we shall describe how to use the Bivariate Beta (see Appendix A) to model the correlation between specificity and sensitivity.

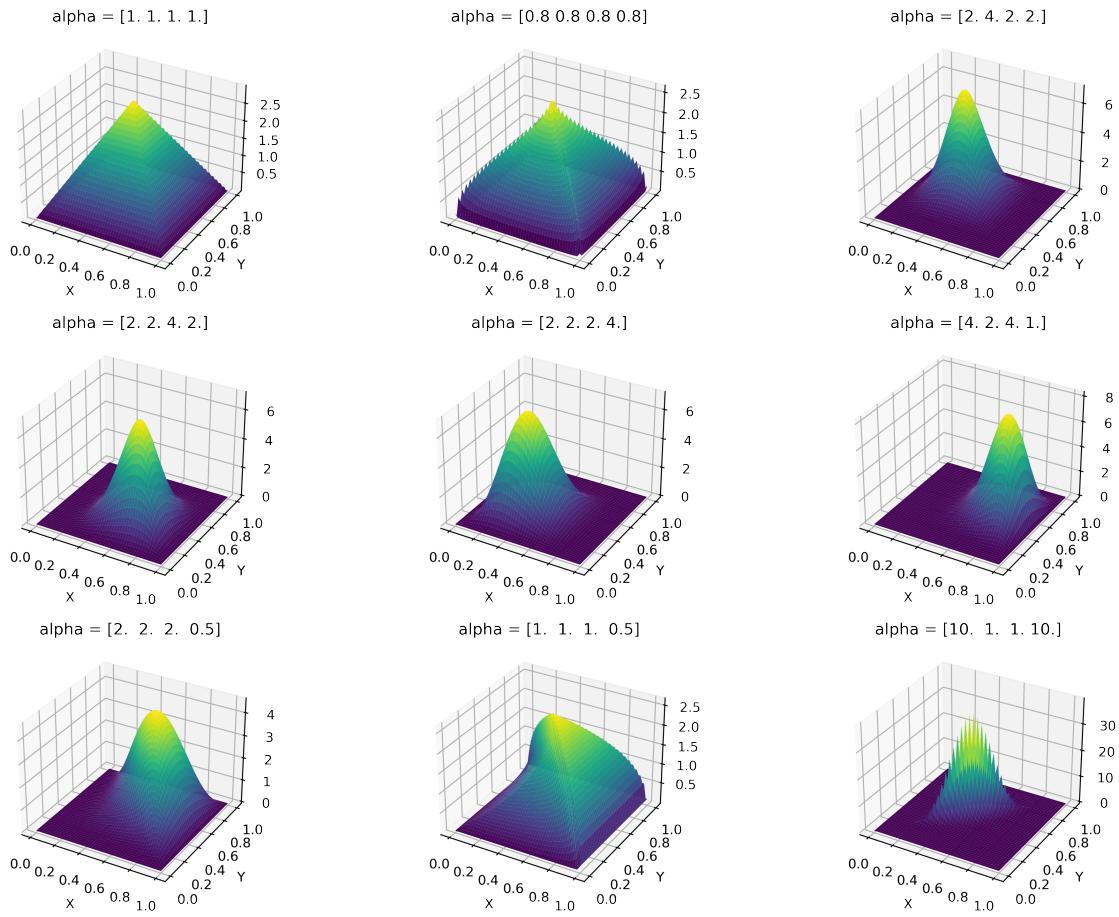


Figure 10 – Different choices of  $\alpha$  and the joint distribution of the variables  $X$  and  $Y$ .

### 3.3 Imperfect tests

This model includes the sensitivity and specificity of the diagnostic test.

$$\begin{aligned}
 T_i &\sim \text{Bernoulli}(p_i) \\
 p_i &= \gamma_s \theta_i + (1 - \gamma_e)(1 - \theta_i), \\
 g(\theta_i) &= g(\theta) + \mathbf{x}_i^T \boldsymbol{\beta}, \\
 \boldsymbol{\beta} &\sim \text{Normal}(\mu, \Sigma), \\
 \theta &\sim \text{Beta}(a^p, b^p) \\
 \gamma_s &\sim \text{Beta}(a^s, b^s), \\
 \gamma_e &\sim \text{Beta}(a^e, b^e),
 \end{aligned} \tag{3.2}$$

where  $a^p, a^s, a^e, b^p, b^s, b^e \in \mathbb{R}_{++}$  are fixed hyperparameters. This model does not include prior knowledge about the correlation between specificity and sensitivity.

### 3.3.1 Simulated data

Consider the following model (GELMAN; CARPENTER, 2020):

$$\begin{aligned} y &\sim \text{Binomial}(n, p), \\ p &= \theta\gamma_s + (1 - \theta)(1 - \gamma_e), \end{aligned}$$

such that  $y$  is the number of positive tests in a population of size  $n$ . In a Bayesian paradigm, a prior  $\pi(\theta, \gamma_e, \gamma_s)$  must be specified. For instance,  $\pi(\theta, \gamma_e, \gamma_s) = \pi(\theta)\pi(\gamma_e, \gamma_s)$  and  $\theta \sim \text{Beta}(\alpha_\theta, \beta_\theta)$ , in which  $\alpha_\theta$  and  $\beta_\theta$  are positive hyperparameters. Since the three parameters  $\theta$ ,  $\gamma_e$ , and  $\gamma_s$  are not jointly identifiable only from  $y$ , prior information on  $\gamma_e$  and  $\gamma_s$  need be added.

## 3.4 Imperfect tests and respondent-driven sampling

For now, we consider the network dependence induced by the RDS with no associated model. Therefore, we treat it as a random effect for each individual. Conditionally autoregressive (CAR) models in the Gaussian case are used. Let  $[\tilde{Q}]_{ij} = \tilde{q}_{ij}$  be a fixed matrix which measures the distance between  $i$  and  $j$ , and  $\tilde{q}_{i+} = \sum_j \tilde{q}_{ij}$ . In general, we use

$$\tilde{q}_{ij} = \begin{cases} 1, & \text{if } i \text{ recruited } j \text{ or the contrary} \\ 0, & \text{otherwise.} \end{cases}$$

Next we define the scaled adjacency matrix  $Q = D^{-1}\tilde{Q}$ , such that  $D$  is a diagonal matrix with  $D_{ii} = \tilde{q}_{i+}$ . Finally let  $|\rho| < 1$  be a parameter to controls the dependence between neighbors. Hence, we specify the model as follows:

$$\begin{aligned} T_i &\sim \text{Bernoulli}(p_i) \\ p_i &= \gamma_s\theta_i + (1 - \gamma_e)(1 - \theta_i), \\ g(\theta_i) &= g(\theta) + \mathbf{x}_i^T \boldsymbol{\beta} + \omega_i, \\ \omega_i | \{\omega_j\}_{j \neq i}, \tau &\sim \text{Normal} \left( \rho \sum_j q_{ij} \omega_j, \tau^{-1} / \tilde{q}_{i+} \right) \\ \boldsymbol{\beta} &\sim \text{Normal}(\mu, \Sigma), \\ \theta &\sim \text{Beta}(a^p, b^p) \\ \gamma_s &\sim \text{Beta}(a^s, b^s), \\ \gamma_e &\sim \text{Beta}(a^e, b^e), \\ \tau &\sim \text{Gamma}(a^\tau, b^\tau). \end{aligned} \tag{3.3}$$

By Brook's Lemma ([BROOK, 1964](#)), the joint distribution of  $\omega$  can be specified as

$$\omega \sim \text{Normal} \left( 0, [\tau(D - \rho\tilde{Q})]^{-1} \right).$$

### 3.4.1 Simulated data

- a) Between the model with the log odds of prevalence having a Gaussian prior distribution and the other with the prevalence having a Beta prior distribution, the latter was usually faster and without divergences. Therefore the preferable model is with the prevalence.
- b) Non-centred distributions are really worst.
- c) Comparison between parametrization of sigma and tau showed that they are similar in sight of time of execution, energy and divergences, among others diagnostics. However, the mean estimate of sigma is more controlled. The median estimate is very similar. This happens because there are a few very high samples for  $\tau$  that will have high weight in the final result. Small samples for  $\sigma$  have less impact, despite having some.
- d) More sparse matrices (RDS data is very sparse) is generating the funnel we do not want to see. This is not connected to the number of connected components. In order to see that, a simple example with the Erdos-Renyi Random Graph can answer to us. In the sparse case, the number of edges is  $O(n)$  with  $p = 1/n$ . If  $p = 1$ , the number of edges is  $O(n^2)$  and the funnel disappears. This problem does not appear in the poisson model.
- e) The effect of  $\rho$  is really observed in the literature in the paper: "A close look at the spatial structure implied by the CAR and SAR models".

### 3.4.2 Exponential Random Graph Model (ERGM)

RDS has the constraint of being without replacement. For that reason, we do not observe all links among the samples ([CRAWFORD, 2016](#)). Considering the model developed by Crawford, we can model the matrix  $Q$  as *Exponential Random Graph Model* (ERGM). Define the following

- a)  $s = \text{tril}(QC)^T \mathbf{1} + C^T u$ , such that  $Q$  is the adjacency matrix of the recruited subjects,  $C$  is the *Coupon Matrix*,  $u$  the vector of the number of edge ends belonging to each vertex (in the order of recruitment) that are not connected to any other sampled vertex, and  $\text{tril}(M)$  the lower triangle of  $M$ .

- b)  $T(Q) = -\lambda \mathbf{s}$ , such that  $\lambda$  is the rate of the recruitment time.
- c)  $V(Q) = \sum_{k \text{ is not seed}} \log(\lambda s_k)$
- d)  $w = (0, t_2 - t_1, \dots, t_n - t_{n-1})$  is thely worst. 3. Comparison between parametrization of sigma and tau showed that they are similar in sight of time of execution, energy and divergences, among others diagnostics. However, the mean estimate of sigma is more controlled. The median estimate is very similar. This happens because there are a few very vector of the waiting times between recruitments.

Therefore  $\Pr(Q|w) \propto \exp[T(Q)^T w + V(Q)]$ . With that, the model becomes

$$\begin{aligned}
 T_i &\sim \text{Bernoulli}(p_i) \\
 p_i &= \gamma_s \theta_i + (1 - \gamma_e)(1 - \theta_i), \\
 g(\theta_i) &= g(\theta) + \mathbf{x}_i^T \boldsymbol{\beta} + \omega_i, \\
 \omega_i | \{\omega_j\}_{j \neq i}, \tau &\sim \text{Normal} \left( \rho \sum_j q_{ij} \omega_j / q_{i+}, \tau^2 / q_{i+} \right) \\
 Q|w &\propto \exp[T(Q)^T w + V(Q)] \\
 \lambda &\sim \Gamma(a^\lambda, b^\lambda), \\
 \boldsymbol{\beta} &\sim \text{Normal}(\mu, \Sigma), \\
 \theta &\sim \text{Beta}(a^p, b^p) \\
 \gamma_s &\sim \text{Beta}(a^s, b^s), \\
 \gamma_e &\sim \text{Beta}(a^e, b^e), \\
 \tau &\sim \text{Normal}^+(0, \sigma_\tau^2).
 \end{aligned} \tag{3.4}$$

The problem with this model is that we are assigning a posterior distribution for  $Q$ .

## 3.5 Model extensions

Several characteristics of RDS were not include in the previous model, such as homophily, bottlenecks, and sampling weights. This section aims to build some options for these aspects and establish future works in that line.

- a) *Homophily model*: (YAUCK et al., 2021)
- b) *Sampling weights*: GLM weighted
- c) *Bottlenecks*

## **3.6 Mispecified data simulation**

## 4 Discussion about prior distributions and sensitivity analysis

**4.1 Prior analysis of sensitivity and specificity**

**4.2 Prior analysis on the parameter tau**

**4.3 Prior analysis on theta**

## 5 Real data applications

## 6 Conclusion

Parte final do trabalho, apresenta as conclusões correspondentes aos objetivos ou hipóteses.

# References

- AVERY, Lisa. **Statistical Methods for Studies Using Respondent Driven Sampling with Applications to Urban Indigenous Health.** 2020. PhD thesis – York University, Toronto, Ontario.
- BANERJEE, Sudipto; CARLIN, Bradley P; GELFAND, Alan E. **Hierarchical modeling and analysis for spatial data.** [S.l.]: Chapman and Hall/CRC, 2003.
- BARAFF, Aaron J; MCCORMICK, Tyler H; RAFTERY, Adrian E. Estimating uncertainty in respondent-driven sampling using a tree bootstrap method. **Proceedings of the National Academy of Sciences**, National Acad Sciences, v. 113, n. 51, p. 14668–14673, 2016.
- BASTOS, Francisco I; BASTOS, Leonardo Soares, et al. HIV, HCV, HBV, and syphilis among transgender women from Brazil: assessing different methods to adjust infection rates of a hard-to-reach, sparse population. **Medicine**, Wolters Kluwer Health, v. 97, 1 Suppl, 2018.
- BASTOS, Leonardo S.; CARVALHO, Luiz M.; GOMES, Marcelo F.C. Modelling misreported data. In: GAMERMAN, Dani et al. **Building a Platform for Data-Driven Pandemic Prediction.** Boca Raton: CRC Press, 2021. chap. 7, p. 113–139.
- BASTOS, Leonardo S.; PINHO, Adriana A., et al. **Binary regression analysis with network structure of respondent-driven sampling data.** [S.l.: s.n.], 2012. arXiv: [1206.5681 \[stat.AP\]](https://arxiv.org/abs/1206.5681).
- BENNETT, Stephen T; STEYVERS, Mark. Estimating COVID-19 antibody seroprevalence in Santa Clara County, California. A re-analysis of Bendavid et al. **MedRxiv**, Cold Spring Harbor Laboratory Press, 2020.
- BETANCOURT, Michael. A conceptual introduction to Hamiltonian Monte Carlo. **arXiv preprint arXiv:1701.02434**, 2017.
- BRANSCUM, AJ; GARDNER, IA; JOHNSON, WO. Estimation of diagnostic-test sensitivity and specificity through Bayesian modeling. **Preventive veterinary medicine**, Elsevier, v. 68, n. 2-4, p. 145–163, 2005.
- BROOK, D. On the distinction between the conditional probability and the joint probability approaches in the specification of nearest-neighbour systems. **Biometrika**, JSTOR, v. 51, n. 3/4, p. 481–483, 1964.

- CARPENTER, Bob et al. Stan: A probabilistic programming language. **Journal of statistical software**, v. 76, n. 1, p. 1–32, 2017.
- CRAWFORD, Forrest W; WU, Jiacheng; HEIMER, Robert. Hidden population size estimation from respondent-driven sampling: a network approach. **Journal of the American Statistical Association**, Taylor & Francis, v. 113, n. 522, p. 755–766, 2018.
- CRAWFORD, Forrest W. The Graphical Structure of Respondent-driven Sampling. **Sociological Methodology**, v. 46, n. 1, p. 187–211, 2016. Available from: <<https://doi.org/10.1177/0081175016641713>>.
- DAMACENA, Giseli Nogueira et al. Application of the Respondent-Driven Sampling methodology in a biological and behavioral surveillance survey among female sex workers, Brazil, 2016. **Revista Brasileira de Epidemiologia**, SciELO Brasil, v. 22, 2019.
- DEAUX, Edward; CALLAGHAN, John W. Key Informant Versus Self-Report Estimates of Health-Risk Behavior. **Evaluation Review**, v. 9, n. 3, p. 365–368, 1985. Available from: <<https://doi.org/10.1177/0193841X8500900308>>.
- FELLOWS, Ian E. Respondent-driven sampling and the homophily configuration graph. **Statistics in medicine**, Wiley Online Library, v. 38, n. 1, p. 131–150, 2019.
- FISHER, Ronald A. On the mathematical foundations of theoretical statistics. **Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character**, The Royal Society London, v. 222, n. 594-604, p. 309–368, 1922.
- GELFAND, Alan E; SAHU, Sujit K. Identifiability, improper priors, and Gibbs sampling for generalized linear models. **Journal of the American Statistical Association**, Taylor & Francis, v. 94, n. 445, p. 247–253, 1999.
- GELMAN, Andrew. Scaling regression inputs by dividing by two standard deviations. **Statistics in medicine**, Wiley Online Library, v. 27, n. 15, p. 2865–2873, 2008.
- GELMAN, Andrew; CARPENTER, Bob. Bayesian analysis of tests with unknown specificity and sensitivity. **Journal of the Royal Statistical Society: Series C (Applied Statistics)**, Wiley Online Library, v. 69, n. 5, p. 1269–1283, 2020.
- GELMAN, Andrew; JAKULIN, Aleks, et al. A weakly informative default prior distribution for logistic and other regression models. **The annals of applied statistics**, Institute of Mathematical Statistics, v. 2, n. 4, p. 1360–1383, 2008.

- GILE, Krista J. Improved inference for respondent-driven sampling data with application to HIV prevalence estimation. **Journal of the American Statistical Association**, Taylor & Francis, v. 106, n. 493, p. 135–146, 2011.
- GILE, Krista J; BEAUDRY, Isabelle S, et al. Methods for inference from respondent-driven sampling data. **Annual Review of Statistics and Its Application**, Annual Reviews, v. 5, p. 65–93, 2018.
- GILE, Krista J; HANDCOCK, Mark S. Network model-assisted inference from respondent-driven sampling data. **Journal of the Royal Statistical Society. Series A,(Statistics in Society)**, NIH Public Access, v. 178, n. 3, p. 619, 2015.
- \_\_\_\_\_. Respondent-driven sampling: An assessment of current methodology. **Sociological methodology**, Sage Publications Sage CA: Los Angeles, CA, v. 40, n. 1, p. 285–327, 2010.
- GILE, Krista J; JOHNSTON, Lisa G; SALGANIK, Matthew J. Diagnostics for respondent-driven sampling. **Journal of the Royal Statistical Society. Series A,(Statistics in Society)**, NIH Public Access, v. 178, n. 1, p. 241, 2015.
- GOEL, Sharad; SALGANIK, Matthew J. Respondent-driven sampling as Markov chain Monte Carlo. **Statistics in medicine**, Wiley Online Library, v. 28, n. 17, p. 2202–2229, 2009.
- GOODMAN, Leo A. Snowball Sampling. **The Annals of Mathematical Statistics**, Institute of Mathematical Statistics, v. 32, n. 1, p. 148–170, 1961.  
Available from: <<https://doi.org/10.1214/aoms/1177705148>>.
- HECKATHORN, Douglas D. Respondent-driven sampling II: deriving valid population estimates from chain-referral samples of hidden populations. **Social problems**, Oxford University Press, v. 49, n. 1, p. 11–34, 2002.
- \_\_\_\_\_. Respondent-Driven Sampling: A New Approach to the Study of Hidden Populations. **Social Problems**, [Oxford University Press, Society for the Study of Social Problems], v. 44, n. 2, p. 174–199, 1997. Available from:  
<<http://www.jstor.org/stable/3096941>>.
- HOFF, Peter D. **A first course in Bayesian statistical methods**. [S.l.]: Springer, 2009. v. 580.
- KÜCHENHOFF, H. The identification of logistic regression models with errors in the variables. **Statistical Papers**, Springer, v. 36, n. 1, p. 41–47, 1995.
- LEHMANN, Eric L. Model specification: the views of Fisher and Neyman, and later developments. In: SELECTED Works of EL Lehmann. [S.l.]: Springer, 2012. P. 955–963.

- LEVIN, David A; PERES, Yuval. **Markov chains and mixing times.** [S.l.]: American Mathematical Soc., 2017. v. 107.
- LIN, Jiayu. On the dirichlet distribution. **Mater's Report**, Queen's University Kingston Ontario, Canada, 2016.
- LINDLEY, Dennis Victor. **Bayesian statistics: A review.** [S.l.]: SIAM, 1972.
- MCCULLAGH, Peter; NELDER, John A. **Generalized linear models.** [S.l.]: Routledge, 2019.
- MCLAUGHLIN, Katherine R. A Bayesian framework for modelling the preferential selection process in respondent-driven sampling. **Statistical Modelling**, SAGE Publications Sage India: New Delhi, India, p. 1471082x211043945, 2021.
- MITCHELL, Stephanie L et al. Performance of SARS-CoV-2 antigen testing in symptomatic and asymptomatic adults: a single-center evaluation. **BMC Infectious Diseases**, Springer, v. 21, n. 1, p. 1–7, 2021.
- MOTA, Rosa Maria Salani. **Respondent driven sampling (RDS) aplicado à população de homens que fazem sexo com homens no Brasil.** 2012. PhD thesis – Universidade Federal do Ceará. Faculdade de Medicina, Fortaleza.
- NOORDZIJ, Marlies et al. Measures of disease frequency: prevalence and incidence. **Nephron Clinical Practice**, Karger Publishers, v. 115, n. 1, p. c17–c20, 2010.
- OGLE, Kiona; BARBER, Jarrett J. Ensuring identifiability in hierarchical mixed effects Bayesian models. **Ecological Applications**, Wiley Online Library, v. 30, n. 7, e02159, 2020.
- OLKIN, Ingram; TRIKALINOS, Thomas A. Constructions for a bivariate beta distribution. **Statistics & Probability Letters**, Elsevier, v. 96, p. 54–60, 2015.
- OTT, Miles Q et al. Reduced bias for respondent-driven sampling: accounting for non-uniform edge sampling probabilities in people who inject drugs in Mauritius. **Journal of the Royal Statistical Society: Series C (Applied Statistics)**, Wiley Online Library, v. 68, n. 5, p. 1411–1429, 2019.
- REITSMA, Johannes B et al. Bivariate analysis of sensitivity and specificity produces informative summary measures in diagnostic reviews. **Journal of Clinical Epidemiology**, Elsevier, v. 58, n. 10, p. 982–990, 2005.
- RIDDELL, Allen; HARTIKAINEN, Ari; CARTER, Matthew. **pystan (3.0.0).** [S.l.: s.n.], Mar. 2021. PyPI.
- ROBERT, Christian. **The Bayesian choice: from decision-theoretic foundations to computational implementation.** [S.l.]: Springer Science & Business Media, 2007.

- ROGAN, Walter J; GLADEN, Beth. Estimating prevalence from the results of a screening test. **American journal of epidemiology**, Oxford University Press, v. 107, n. 1, p. 71–76, 1978.
- ROTHMAN, Kenneth J; GREENLAND, Sander; LASH, Timothy L, et al. **Modern epidemiology**. [S.l.]: Wolters Kluwer Health/Lippincott Williams & Wilkins Philadelphia, 2008. v. 3.
- RUTJES, AWS et al. Evaluation of diagnostic tests when there is no gold standard. A review of methods. **HEALTH TECHNOLOGY ASSESSMENT-SOUTHAMPTON-**, National Coordinating Centre for Health Technology Assessment, v. 11, n. 50, 2007.
- SALGANIK, Matthew J. Variance estimation, design effects, and sample size calculations for respondent-driven sampling. **Journal of Urban Health**, Springer, v. 83, n. 1, p. 98, 2006.
- SALGANIK, Matthew J; FAZITO, Dimitri, et al. Assessing network scale-up estimates for groups most at risk of HIV/AIDS: evidence from a multiple-method study of heavy drug users in Curitiba, Brazil. **American journal of epidemiology**, Oxford University Press, v. 174, n. 10, p. 1190–1196, 2011.
- SALGANIK, Matthew J; HECKATHORN, Douglas D. Sampling and estimation in hidden populations using respondent-driven sampling. **Sociological methodology**, Wiley Online Library, v. 34, n. 1, p. 193–240, 2004.
- SCHERVISH, Mark J. **Theory of statistics**. [S.l.]: Springer Science & Business Media, 2012.
- ŠIMUNDIĆ, Ana-Maria. Measures of diagnostic accuracy: basic definitions. **Ejifcc**, International Federation of Clinical Chemistry and Laboratory Medicine, v. 19, n. 4, p. 203, 2009.
- SPILLER, Michael. **Regression modeling of data collected using respondentdriven sampling**. 2009. PhD thesis – Cornell University.
- STATISTICAT, LLC. LaplacesDemon: A Complete Environment for Bayesian Inference within R. **R Package version**, v. 17, p. 2016, 2016.
- TOLEDO, Lidiane et al. Putting respondent-driven sampling on the map: insights from Rio de Janeiro, Brazil. **JAIDS Journal of Acquired Immune Deficiency Syndromes**, LWW, v. 57, s136–s143, 2011.
- VERSI, E. " Gold standard" is an appropriate term. **BMJ: British Medical Journal**, BMJ Publishing Group, v. 305, n. 6846, p. 187, 1992.

- VOLZ, Erik; HECKATHORN, Douglas D. Probability based estimation theory for respondent driven sampling. **Journal of Official Statistics**, Statistics Sweden (SCB), v. 24, n. 1, p. 79, 2008.
- WATTERS, John K.; BIERNACKI, Patrick. Targeted Sampling: Options for the Study of Hidden Populations. **Social Problems**, Oxford University Press, Society for the Study of Social Problems, v. 36, n. 4, p. 416–430, 1989. Available from: <<http://www.jstor.org/stable/800824>>.
- WORLD HEALTH ORGANIZATION. **Introduction to HIV/AIDS and sexually transmitted infection surveillance: Module 4: Introduction to respondent-driven sampling**. [S.l.], 2013. 389 p., 30 cm. Available from: <<https://apps.who.int/iris/handle/10665/116864>>.
- XIE, Yang; CARLIN, Bradley P. Measures of Bayesian learning and identifiability in hierarchical models. **Journal of Statistical Planning and Inference**, Elsevier, v. 136, n. 10, p. 3458–3477, 2006.
- YAUCK, Mamadou et al. General regression methods for respondent-driven sampling data. **Statistical methods in medical research**, SAGE Publications Sage UK: London, England, v. 30, n. 9, p. 2105–2118, 2021.

# **Appendix**

# APPENDIX A – Bivariate Beta distribution

Let  $U = (U_1, U_2, U_3, U_4) \sim \text{Dirichlet}(\boldsymbol{\alpha})$ , where  $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \alpha_3, \alpha_4)$  with  $\alpha_i > 0, i = 1, \dots, 4$  and  $U_4 = 1 - U_1 + U_2 + U_3$ . The joint density of  $U$  with respect to the Lebesgue measure is given by

$$f_U(u_1, u_2, u_3) = \frac{1}{B(\boldsymbol{\alpha})} u_1^{\alpha_1-1} u_2^{\alpha_2-1} u_3^{\alpha_3-1} (1 - u_1 - u_2 - u_3)^{\alpha_4-1}, \quad (\text{A.1})$$

when  $u_i \in [0, 1], i = 1, 2, 3$ ,  $u_1 + u_2 + u_3 \leq 1$ , and 0 otherwise. The normalizing constant is, for  $v \in \mathbb{R}^n$ ,

$$B(v) = \frac{\prod_{i=1}^n \Gamma(v_i)}{\Gamma(\sum_{i=1}^n v_i)}.$$

**Definition A.0.1.** Let

$$X = U_1 + U_2 \text{ and } Y = U_1 + U_3. \quad (\text{A.2})$$

The distribution of  $(X, Y)$  is *Bivariate Beta* with parameters  $\boldsymbol{\alpha}$ .

**Proposition A.0.1.** *The marginal distribution of  $X$  is Beta with parameters  $\alpha_1 + \alpha_2$  and  $\alpha_3 + \alpha_4$ . Similarly, the marginal distribution of  $Y$  is Beta with parameters  $\alpha_1 + \alpha_3$  and  $\alpha_2 + \alpha_4$ .*

*Proof.* First we derive the probability density of  $(U_1, U_2)$  with respect to the Lebesgue measure.

$$\begin{aligned} f_{U_1, U_2}(u_1, u_2) &= \int_{-\infty}^{\infty} f_U(u_1, u_2, u_3) du_3 \\ &= \frac{1}{B(\boldsymbol{\alpha})} \int_0^1 u_1^{\alpha_1-1} u_2^{\alpha_2-1} u_3^{\alpha_3-1} (1 - u_1 - u_2 - u_3)^{\alpha_4-1} du_3 \\ &= \frac{1}{B(\boldsymbol{\alpha})} u_1^{\alpha_1-1} u_2^{\alpha_2-1} \int_0^1 u_3^{\alpha_3-1} (1 - u_1 - u_2 - u_3)^{\alpha_4-1} du_3. \end{aligned} \quad (\text{A.3})$$

Let  $u_3 = (1 - u_1 - u_2)z$ . Then,

$$\begin{aligned} f_{U_1, U_2}(u_1, u_2) &= \frac{1}{B(\boldsymbol{\alpha})} u_1^{\alpha_1-1} u_2^{\alpha_2-1} \int_0^1 (1 - u_1 - u_2)^{\alpha_3-1} z^{\alpha_3-1} (1 - u_1 - u_2)^{\alpha_4-1} (1 - z)^{\alpha_4-1} dz. \\ &= \frac{1}{B(\boldsymbol{\alpha})} u_1^{\alpha_1-1} u_2^{\alpha_2-1} (1 - u_1 - u_2)^{\alpha_3+\alpha_4-1} \int_0^1 z^{\alpha_3-1} (1 - z)^{\alpha_4-1} dz. \\ &= \frac{1}{B(\boldsymbol{\alpha})} u_1^{\alpha_1-1} u_2^{\alpha_2-1} (1 - u_1 - u_2)^{\alpha_3+\alpha_4-1} \frac{\Gamma(\alpha_3)\Gamma(\alpha_4)}{\Gamma(\alpha_3 + \alpha_4)} \\ &= \frac{1}{B(\alpha_1, \alpha_2, \alpha_3 + \alpha_4)} u_1^{\alpha_1-1} u_2^{\alpha_2-1} (1 - u_1 - u_2)^{\alpha_3+\alpha_4-1}. \end{aligned} \quad (\text{A.4})$$

We conclude that

$$(U_1, U_2, 1 - U_1 - U_2) \sim \text{Dirichlet}(\alpha_1, \alpha_2, \alpha_3 + \alpha_4).$$

Define

$$H(v) = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} v, \text{ for } v \in \mathbb{R}^2.$$

Then  $(U_1, X) = H(U_1, U_2)$  and  $H(\cdot)$  is bijective and differentiable function.

By the Change of Variable Formula,

$$\begin{aligned} f_{U_1, X}(u_1, x) &= f(H^{-1}(u_1, x)) \left| \det \left[ \frac{dH^{-1}(v)}{dv} \right]_{v=(u_1, x)} \right| \\ &= f(u_1, x - u_1) = \frac{1}{B(\alpha_1, \alpha_2, \alpha_3 + \alpha_4)} u_1^{\alpha_1-1} (x - u_1)^{\alpha_2-1} (1 - x)^{\alpha_3 + \alpha_4 - 1}, \end{aligned} \quad (\text{A.5})$$

where  $(u_1, x)$  belongs to the triangle defined by the points  $(0,0)$ ,  $(0,1)$ , and  $(1,1)$ .

The distribution of  $X$  for  $x \in [0, 1]$  is

$$\begin{aligned} f_X(x) &= \frac{1}{B(\alpha_1, \alpha_2, \alpha_3 + \alpha_4)} \int_0^x u_1^{\alpha_1-1} (x - u_1)^{\alpha_2-1} (1 - x)^{\alpha_3 + \alpha_4 - 1} du_1 \\ &= \frac{1}{B(\alpha_1, \alpha_2, \alpha_3 + \alpha_4)} (1 - x)^{\alpha_3 + \alpha_4 - 1} \int_0^x u_1^{\alpha_1-1} (x - u_1)^{\alpha_2-1} du_1. \\ &= \frac{1}{B(\alpha_1, \alpha_2, \alpha_3 + \alpha_4)} (1 - x)^{\alpha_3 + \alpha_4 - 1} \int_0^x x^{\alpha_1-1} \left( \frac{u_1}{x} \right)^{\alpha_1-1} x^{\alpha_2-1} \left( 1 - \frac{u_1}{x} \right)^{\alpha_2-1} du_1. \end{aligned} \quad (\text{A.6})$$

Setting  $u = u_1/x$  (if  $x = 0$ ,  $f_X(x) = 0$ , then suppose  $x > 0$ ), we have,

$$\begin{aligned} f_X(x) &= \frac{1}{B(\alpha_1, \alpha_2, \alpha_3 + \alpha_4)} (1 - x)^{\alpha_3 + \alpha_4 - 1} x^{\alpha_1 + \alpha_2 - 1} \int_0^1 u^{\alpha_1-1} (1 - u)^{\alpha_2-1} du. \\ &= \frac{1}{B(\alpha_1, \alpha_2, \alpha_3 + \alpha_4)} (1 - x)^{\alpha_3 + \alpha_4 - 1} x^{\alpha_1 + \alpha_2 - 1} B(\alpha_1, \alpha_2) \\ &= \frac{1}{B(\alpha_1 + \alpha_2, \alpha_3 + \alpha_4)} (1 - x)^{\alpha_3 + \alpha_4 - 1} x^{\alpha_1 + \alpha_2 - 1} \end{aligned} \quad (\text{A.7})$$

Therefore  $X \sim \text{Beta}(\alpha_1 + \alpha_2, \alpha_3 + \alpha_4)$ . Similarly  $Y \sim \text{Beta}(\alpha_1 + \alpha_3, \alpha_2 + \alpha_4)$ .

□

**Proposition A.0.2.** *The joint density of  $(X, Y)$  with respect to the Lebesgue measure is given by*

$$f_{X,Y}(x, y) = \frac{1}{B(\boldsymbol{\alpha})} \int_{\Omega} u_1^{\alpha_1-1} (x - u_1)^{\alpha_2-1} (y - u_1)^{\alpha_3-1} (1 - x - y + u_1)^{\alpha_4-1} du_1, \quad (\text{A.8})$$

where

$$\Omega = (\max(0, x + y - 1), \min(x, y)).$$

*Proof.* Note that

$$\begin{bmatrix} U_1 \\ X \\ Y \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} U_1 \\ U_2 \\ U_3 \end{bmatrix},$$

where the linear function is bijective and differentiable function, such that the determinant of the derivative is 1. By the Change of Variable Formula,

$$\begin{aligned} f_{U_1, X, Y}(u_1, x, y) &= f_{U_1, U_2, U_3}(u_1, x - u_1, y - u_2) \\ &= \frac{1}{B(\boldsymbol{\alpha})} u_1^{\alpha_1-1} (x - u_1)^{\alpha_2-1} (y - u_1)^{\alpha_3-1} (1 - x - y + u_1)^{\alpha_4-1}, \end{aligned} \quad (\text{A.9})$$

where  $0 \leq u_1 \leq x, u_1 \leq y$ , and  $0 \leq 1 - x - y + u_1$ . Hence,

$$f_{X, Y}(x, y) = \frac{1}{B(\boldsymbol{\alpha})} \int_{\Omega} u_1^{\alpha_1-1} (x - u_1)^{\alpha_2-1} (y - u_1)^{\alpha_3-1} (1 - x - y + u_1)^{\alpha_4-1} du_1, \quad (\text{A.10})$$

such that  $\Omega = \{u_1 : \max(0, x + y - 1) < u_1 < \min(x, y)\}$ .  $\square$

**Proposition A.0.3.** *The covariance between  $X$  and  $Y$  is*

$$\text{Cov}(X, Y) = \frac{1}{\tilde{\alpha}^2(\tilde{\alpha} + 1)} (\alpha_1\alpha_4 - \alpha_2\alpha_3).$$

*Proof.* Let  $\tilde{\alpha} = \sum_i \alpha_i$ . The covariance between  $U_i$  and  $U_j$  is (LIN, 2016)

$$\text{Cov}(U_i, U_j) = -\frac{\alpha_i\alpha_j}{\tilde{\alpha}^2(\tilde{\alpha} + 1)}, i, j = 1, \dots, 4, i \neq j \quad (\text{A.11})$$

and the variance of  $U_i$  is

$$\text{Var}(U_i) = \frac{\alpha_i(\tilde{\alpha} - \alpha_i)}{\tilde{\alpha}^2(\tilde{\alpha} + 1)}, \quad (\text{A.12})$$

since  $U_i \sim \text{Beta}(\alpha_i, \tilde{\alpha} - \alpha_i)$ . Therefore

$$\text{Cov}(X, Y) = \text{Cov}(U_1 + U_2, U_1 + U_3) = \frac{1}{\tilde{\alpha}^2(\tilde{\alpha} + 1)} (\alpha_1\alpha_4 - \alpha_2\alpha_3) \quad (\text{A.13})$$

$\square$

The main moments of  $X$  and  $Y$  are the following

$$\begin{aligned} \mathbb{E}(X) &= \mathbb{E}(U_1 + U_2) = \frac{\alpha_1 + \alpha_2}{\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4} \\ \mathbb{E}(Y) &= \mathbb{E}(U_1 + U_3) = \frac{\alpha_1 + \alpha_3}{\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4} \\ \text{Var}(X) &= \text{Cov}(U_1 + U_2, U_1 + U_2) = \frac{1}{\tilde{\alpha}^2(\tilde{\alpha} + 1)} (\alpha_1 + \alpha_2)(\alpha_3 + \alpha_4) \\ \text{Var}(Y) &= \text{Cov}(U_1 + U_3, U_1 + U_3) = \frac{1}{\tilde{\alpha}^2(\tilde{\alpha} + 1)} (\alpha_1 + \alpha_3)(\alpha_2 + \alpha_4) \\ \text{Cor}(X, Y) &= \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}} = \frac{\alpha_1\alpha_4 - \alpha_2\alpha_3}{\sqrt{(\alpha_1 + \alpha_2)(\alpha_3 + \alpha_4)(\alpha_1 + \alpha_3)(\alpha_2 + \alpha_4)}} \end{aligned}$$

*Remark A.0.1.* The original paper has a mistake at page 6.

## A.1 Comments about integration

The density of  $(X, Y)$  with respect to the Lebesgue measure is  $f_{X,Y}(x, y)$  as in equation (A.10). Therefore it can be undefined in sets of null Lebesgue measure in  $\mathbb{R}^2$ . This section aims to find them to help writing the function properly. If  $\alpha_i \geq 1$ ,  $i = 1, \dots, 4$ , the integral is clearly well defined for every  $x, y \in [0, 1]$ . Let  $0 < \alpha_2 = \alpha_3 = a \leq 0.5$  and  $x = y < 0.5$ . Then

$$\begin{aligned} f_{X,Y}(x, y) &= \frac{1}{B(\boldsymbol{\alpha})} \int_0^x u_1^{\alpha_1-1} (x - u_1)^{a-1} (x - u_1)^{a-1} (1 - 2x + u_1)^{\alpha_4-1} du_1 \\ &= \frac{1}{B(\boldsymbol{\alpha})} \int_0^{x/2} u_1^{\alpha_1-1} (x - u_1)^{2a-2} (1 - 2x + u_1)^{\alpha_4-1} du_1 + \\ &\quad + \frac{1}{B(\boldsymbol{\alpha})} \int_{x/2}^x u_1^{\alpha_1-1} (x - u_1)^{2a-2} (1 - 2x + u_1)^{\alpha_4-1} du_1 \end{aligned}$$

Note that the first integral is well defined and non-negative. If  $\alpha_1 \geq 1$ ,

$$\begin{aligned} &\int_0^{x/2} u_1^{\alpha_1-1} (x - u_1)^{2a-2} (1 - 2x + u_1)^{\alpha_4-1} du_1 \\ &\leq \int_0^{x/2} \frac{x^{\alpha_1-1}}{2} \left(\frac{x}{2}\right)^{2a-2} \max\left(\left(1 - \frac{3}{2}x\right)^{\alpha_4-1}, (1 - 2x)^{\alpha_4-1}\right) du_1 < +\infty. \end{aligned}$$

If  $0 < \alpha_1 < 1$ ,

$$\begin{aligned} &\int_0^{x/2} u_1^{\alpha_1-1} (x - u_1)^{2a-2} (1 - 2x + u_1)^{\alpha_4-1} du_1 \\ &= \lim_{t \rightarrow 0^+} \int_t^{x/2} u_1^{\alpha_1-1} \left(\frac{x}{2}\right)^{2a-2} \max\left(\left(1 - \frac{3}{2}x\right)^{\alpha_4-1}, (1 - 2x)^{\alpha_4-1}\right) du_1 \\ &= K(x) \lim_{t \rightarrow 0^+} \int_t^{x/2} u_1^{\alpha_1-1} du_1 \\ &= \frac{K(x)}{\alpha_1} \lim_{t \rightarrow 0^+} \left[\left(\frac{x}{2}\right)^{\alpha_1} - t^{\alpha_1}\right] < +\infty. \end{aligned}$$

where  $K(x)$  is a function of  $x$ . Moreover, since the integrand is non-negative, so is the integral. On the other hand, the second integral is not defined:

$$\begin{aligned} &\int_{x/2}^x u_1^{\alpha_1-1} (x - u_1)^{2a-2} (1 - 2x + u_1)^{\alpha_4-1} du_1 \\ &\geq \int_{x/2}^x \min\left(\left(\frac{x}{2}\right)^{\alpha_1-1}, x^{\alpha_1-1}\right) (x - u_1)^{2a-2} \min\left(\left(1 - \frac{3}{2}x\right)^{\alpha_4-1}, (1 - x)^{\alpha_4-1}\right) du_1 \\ &= K'(x) \int_0^{x/2} v^{2a-2} dv \\ &= \begin{cases} \frac{K'(x)}{2a-1} \lim_{t \rightarrow 0^+} [(x/2)^{2a-1} - t^{2a-1}] & \text{if } a < 0.5 \\ K'(x) \lim_{t \rightarrow 0^+} [\log(x/2) - \log(t)] & \text{if } a = 0.5 \end{cases} \\ &\rightarrow +\infty. \end{aligned}$$

Based on this divergence, we conclude that if  $0 < \alpha_2 = \alpha_3 \leq 0.5$  and  $x = y < 0.5$ ,  $f_{X,Y}(x, y)$  is not defined. Note that if  $x = y \geq 0.5$ , divergence problems still happens, since the problems appear when  $u_1$  converges to  $x$ . Similar calculations show that if  $x + y = 1$  and  $0 < \alpha_1 = \alpha_4 \leq 0.5$ , the density is also not defined. More generally,  $f_{X,Y}(x, y)$  is not defined if

1.  $\alpha_1 + \alpha_4 \leq 1$  and  $x + y = 1$ .
2.  $\alpha_2 + \alpha_3 \leq 1$  and  $x = y$ .

## A.2 Specifying parameters $\alpha$

Suppose that the researcher has knowledge about the main moments of  $X$  and  $Y$ , such that  $\mathbb{E}(X) = m_1 \in (0, 1)$ ,  $\mathbb{E}(Y) = m_2 \in (0, 1)$ ,  $\text{Var}(X) = v_1 \in (0, 1)$ , and  $\text{Var}(Y) = v_2 \in (0, 1)$ . Notice that  $v_1 + m_1^2 = \text{Var}(X_1) + \mathbb{E}[X_1]^2 = \mathbb{E}[X_1^2]$  and

$$\mathbb{E}[X_1^2] - \mathbb{E}[X_1] = \frac{(\alpha_1 + \alpha_2 + 1)(\alpha_1 + \alpha_2)}{(\tilde{\alpha} + 1)\tilde{\alpha}} - \frac{\alpha_1 + \alpha_2}{\tilde{\alpha}} = -\frac{(\alpha_1 + \alpha_2)(\alpha_3 + \alpha_4)}{\tilde{\alpha}(\tilde{\alpha} + 1)} < 0,$$

that is,  $v_1 + m_1^2 - m_1 < 0 \implies v_1 < m_1 - m_1^2$  and similarly,  $v_2 < m_2 - m_2^2$ . After fixing these quantities, we will have a non-linear system with four equations and four unknown variables. Hence, we want to solve the following

$$\begin{cases} m_1 = \frac{\alpha_1 + \alpha_2}{\tilde{\alpha}} \\ m_2 = \frac{\alpha_1 + \alpha_3}{\tilde{\alpha}} \\ v_1 = \frac{(\alpha_1 + \alpha_2)(\alpha_3 + \alpha_4)}{\tilde{\alpha}^2(\tilde{\alpha} + 1)} = m_1 \frac{\alpha_3 + \alpha_4}{\tilde{\alpha}(\tilde{\alpha} + 1)} \\ v_2 = \frac{(\alpha_1 + \alpha_3)(\alpha_2 + \alpha_4)}{\tilde{\alpha}^2(\tilde{\alpha} + 1)} = m_2 \frac{\alpha_2 + \alpha_4}{\tilde{\alpha}(\tilde{\alpha} + 1)}. \end{cases} \quad (\text{A.14})$$

**Proposition A.2.1.** *System (A.14) has a solution if, and only if, the relation*

$$v_2 = \frac{(1 - m_2)\tilde{\alpha}}{\tilde{\alpha}(\tilde{\alpha} + 1)} = \frac{1 - m_2}{\frac{m_1 - m_1^2}{v_1}} = \frac{v_1(1 - m_2)}{m_1(1 - m_1)}, \quad (\text{A.15})$$

*is satisfied. When there is a solution, there will be infinitely many and they all lay in the ray*

$$\mathcal{L} = \{(1, -1, -1, 1)\alpha_4 + k : \alpha_4 > 0\},$$

*such that  $k = ((m_1 + m_2 - 1)\tilde{\alpha}, (1 - m_2)\tilde{\alpha}, (1 - m_1)\tilde{\alpha}, 0)$ .*

*Proof.* The first two equations of the system (A.14) can be rewritten as a linear system:

$$\begin{aligned}(m_1 - 1)\alpha_1 + (m_1 - 1)\alpha_2 + m_1\alpha_3 + m_1\alpha_4 &= 0 \\ (m_2 - 1)\alpha_1 + m_2\alpha_2 + (m_2 - 1)\alpha_3 + m_2\alpha_4 &= 0,\end{aligned}$$

which is equivalent to

$$\begin{aligned}\alpha_1 + \alpha_2 + \frac{m_1}{m_1 - 1}\alpha_3 + \frac{m_1}{m_1 - 1}\alpha_4 &= 0 \\ \alpha_2 + \frac{1 - m_2}{m_1 - 1}\alpha_3 + \frac{m_1 - m_2}{m_1 - 1}\alpha_4 &= 0.\end{aligned}$$

Then, we can write  $\alpha_1$  and  $\alpha_2$  as functions of  $\alpha_3$  and  $\alpha_4$ :

$$\alpha_1 = \frac{m_1 + m_2 - 1}{1 - m_1}\alpha_3 + \frac{m_2}{1 - m_1}\alpha_4 \quad (\text{A.16})$$

$$\alpha_2 = \frac{1 - m_2}{1 - m_1}\alpha_3 + \frac{m_1 - m_2}{1 - m_1}\alpha_4. \quad (\text{A.17})$$

With that expression, let  $\alpha_1 = a_3\alpha_3 + a_4\alpha_4$  and  $\alpha_2 = b_3\alpha_3 + b_4\alpha_4$ . Denote  $c_3 = a_3 + b_3 + 1$  and  $c_4 = a_4 + b_4 + 1$ . Then, consider the third equation of the system (A.14),

$$\begin{aligned}\frac{v_1}{m_1} = \frac{\alpha_3 + \alpha_4}{\tilde{\alpha}(\tilde{\alpha} + 1)} &= \frac{\alpha_3 + \alpha_4}{(\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4)^2 + (\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4)} \\ \implies \frac{v_1}{m_1}(\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4)^2 &= \alpha_3 + \alpha_4 - \frac{v_1}{m_1}(\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4) \\ \implies \frac{v_1}{m_1}(c_3\alpha_3 + c_4\alpha_4)^2 &= \left(1 - \frac{v_1}{m_1}c_3\right)\alpha_3 + \left(1 - \frac{v_1}{m_1}c_4\right)\alpha_4 \\ \implies \frac{v_1c_3^2}{m_1}\alpha_3^2 + \left(\frac{2v_1c_3c_4\alpha_4 + v_1c_3}{m_1} - 1\right)\alpha_3 + \left(\frac{v_1c_4^2\alpha_4^2 + v_1c_4\alpha_4}{m_1} - \alpha_4\right) &= 0 \\ \implies v_1c_3^2\alpha_3^2 + (2v_1c_3c_4\alpha_4 + v_1c_3 - m_1)\alpha_3 + (v_1c_4^2\alpha_4^2 + v_1c_4\alpha_4 - m_1\alpha_4) &= 0.\end{aligned}$$

Using a Computer Algebra System (CAS) with the Python library SymPy, the above expression can be simplified as follows:

$$v_1\alpha_3^2 + (v_1(1 - m_1) + 2v_1\alpha_4 - m_1(1 - m_1)^2)\alpha_3 - \alpha_4m_1(1 - m_1)^2 + \alpha_4v_1(1 - m_1) + v_1\alpha_4^2 = 0.$$

This way, the solutions of the above equation are function of  $\alpha_4$ . Therefore, after solving the equations, we can use the last equation of the system (A.14) as a function on of  $\alpha_4$ . Let,

$$\Lambda = (v_1(1 - m_1) + v_1\alpha_4 - m_1(1 - m_1)^2).$$

Then,

$$\begin{aligned}
\Delta &= \left( v_1(1 - m_1) + 2v_1\alpha_4 - m_1(1 - m_1)^2 \right)^2 - 4v_1(\alpha_4 v_1(1 - m_1) - \alpha_4 m_1(1 - m_1)^2 + v_1\alpha_4^2), \\
&= (\Lambda + v_1\alpha_4)^2 - 4v_1\alpha_4\Lambda \\
&= \Lambda^2 - 2\Lambda v_1\alpha_4 + (v_1\alpha_4)^2 \\
&= (\Lambda - v_1\alpha_4)^2 \\
&= \left( v_1(1 - m_1) - m_1(1 - m_1)^2 \right)^2 \\
&= (1 - m_1)^2(v_1 + m_1^2 - m_1)^2.
\end{aligned}$$

Note that  $v_1 + m_1^2 = \text{Var}(X_1) + \mathbb{E}[X_1]^2 = \mathbb{E}[X_1^2]$  and

$$\mathbb{E}[X_1^2] - \mathbb{E}[X_1] = \frac{(\alpha_1 + \alpha_2 + 1)(\alpha_1 + \alpha_2)}{(\tilde{\alpha} + 1)\tilde{\alpha}} - \frac{\alpha_1 + \alpha_2}{\tilde{\alpha}} = -\frac{(\alpha_1 + \alpha_2)(\alpha_3 + \alpha_4)}{\tilde{\alpha}(\tilde{\alpha} + 1)} < 0.$$

Therefore,

$$\sqrt{\Delta} = (1 - m_1)(m_1 - v_1 - m_1^2)$$

and

$$\begin{aligned}
\alpha_3 &= \frac{1}{2v_1} \left( (m_1(1 - m_1)^2 - v_1(1 - m_1) - 2v_1\alpha_4) \pm (1 - m_1)(m_1 - v_1 - m_1^2) \right) \\
&= -\alpha_4 + \frac{(1 - m_1)(m_1 - m_1^2 - v_1) \pm (1 - m_1)(m_1 - v_1 - m_1^2)}{2v_1}.
\end{aligned}$$

When the sign is negative, we have that  $\alpha_3 = -\alpha_4$ , an impossible solution. Then,

$$\alpha_3 = \frac{(1 - m_1)(m_1 - m_1^2 - v_1)}{v_1} - \alpha_4.$$

We summarize the expressions in function of  $\alpha_4$ :

$$\begin{aligned}
\alpha_3 &= \frac{(1 - m_1)(m_1 - m_1^2 - v_1)}{v_1} - \alpha_4 \\
\alpha_1 &= \frac{m_1 + m_2 - 1}{1 - m_1}\alpha_3 + \frac{m_2}{1 - m_1}\alpha_4 = \frac{(m_1 + m_2 - 1)(m_1 - m_1^2 - v_1)}{v_1} + \alpha_4 \\
\alpha_2 &= \frac{1 - m_2}{1 - m_1}\alpha_3 + \frac{m_1 - m_2}{1 - m_1}\alpha_4 = \frac{(1 - m_2)(m_1 - m_1^2 - v_1)}{v_1} - \alpha_4.
\end{aligned}$$

From here, one can calculate that

$$\tilde{\alpha} = \frac{m_1 - m_1^2 - v_1}{v_1}.$$

Since  $\alpha_2 + \alpha_4 = (1 - m_2)\tilde{\alpha}$ , we have that the last equation of the system (A.14) is given by (A.15), that is, the system (A.14) has a solution if and only if, equation (A.15) is satisfied. If it is, the solution is the ray

$$\mathcal{L} = \{(1, -1, -1, 1)\alpha_4 + k : \alpha_4 > 0\},$$

such that  $k = ((m_1 + m_2 - 1)\tilde{\alpha}, (1 - m_2)\tilde{\alpha}, (1 - m_1)\tilde{\alpha}, 0)$ .

□

Now change the fourth equation of (A.14) by:

$$\text{Cor}(X, Y) = \frac{\alpha_1\alpha_4 - \alpha_2\alpha_3}{\sqrt{(\alpha_1 + \alpha_2)(\alpha_3 + \alpha_4)(\alpha_1 + \alpha_3)(\alpha_2 + \alpha_4)}} = \frac{\alpha_1\alpha_4 - \alpha_2\alpha_3}{\tilde{\alpha}^2\sqrt{m_1m_2(1 - m_1)(1 - m_2)}}$$

Supposing the expression for  $\alpha_1, \alpha_2$  and  $\alpha_3$ , that is,  $m_1, m_2$  and  $v_1$  are fixed, and supposing we fix  $\rho = \text{Cor}(X, Y)$ , we can simplify the above expression (using a software) as follows:

$$\rho = \frac{1}{\tilde{\alpha}\sqrt{m_1m_2(1 - m_1)(1 - m_2)}}\alpha_4 - \sqrt{\frac{(1 - m_1)(1 - m_2)}{m_1m_2}},$$

which is linear on  $\alpha_4$ , that is, for fixed values of  $m_1, m_2, v_1$  and  $\rho$ , there is an unique  $\alpha_4$ , and hence,  $\alpha_1, \alpha_2$  and  $\alpha_3$  that satisfies system (A.14) with the fourth equation changed by the correlation.

## APPENDIX B – Stan codes