

POLYTECH NICE SOPHIA ANTIPOLIS

MASTER IFI/ UBINET TRACK

FINAL REPORT

Network sampling and discovery processes

Author:

Alina TUHOLUKOVA

Supervisors:

Konstantin AVRACHENKOV

Giovanni NEGLIA

August 24, 2015



Contents

1	Respondent-driven sampling	5
1.1	Motivation	5
1.2	Technique of respondent-driven sampling	6
1.3	Current estimators	6
1.3.1	Sample average	7
1.3.2	Volz-Heckathorn estimator	7
1.3.3	Group estimator	7
1.4	Problems of RDS	8
1.4.1	RDS is not uniform	8
1.4.2	RDS is not independent	9
1.4.3	Other problems	10
1.5	Enhanced RDS	10
2	Mathematical model	14
2.1	Network modeling	14
2.2	Network with values	14
2.2.1	Motivation	14
2.2.2	Definitions	15
2.2.3	Gibbs distribution	16
2.2.4	Algorithm	16
2.2.5	Explanatory example	17
2.2.6	Demonstration of random graphs with values	18
2.3	Expected energy in steady state	20
2.3.1	Motivation	20
2.3.2	Analysis	21
2.3.3	Results	25
2.4	Error prediction	26
2.4.1	Variance prediction	26
2.4.2	Geometric correlation	27
2.4.3	Variance with skipping	29
2.4.4	In RDS context	30
2.4.5	Trying to use the result	32
2.4.6	General case	32
2.4.7	Error prediction	34
2.4.8	Discussion	35
2.5	Data	35

2.5.1	Data from the Project 90	35
2.5.2	Data from the Add health project	35
2.5.3	Simulated form Gibbs distribution	36
2.5.4	Experiments	36
3	Comparison to other methods	39
3.1	Further study	39

Introduction

We are living in the era of information when it is crucial to collect data, to be able to analyze them and draw potentially valuable conclusions. Particularly it is interesting to analyze network structures. The examples are: online-social networks, peer-to-peer networks, real social networks, hidden populations.

Online-social network are thriving nowadays. The most popular are: Google+ (around 1.6 billion users), Facebook (around 1.28 billion users), Twitter (around 645 million users), Instagram (around 300 million users), VK (around 250 million users), LinkedIn (around 200 million users). These networks gather a lot of valuable information like users interests, users characteristics, etc. Great part of it is in the free access. This information can facilitate the life to the sociologists and give them another instrument for the research.

As well as it interesting to analyze the structure of network graph by itself it can be useful to study the connection between network structure and the characteristics of the users. It is frequently observed fact that friends tend to be similar interests. It can be the influence of your friend that you listen to the rock music or the opposite you both are fond of rock music and therefore you are friends. One way or another this property of sharing the common characteristics between people that are in contact is usually observed in the networks and is called *homophily*.

And actually sometimes it is the common interest that is in foundation of the network. Online social network Flixster can help to meet people with the similar tastes in movies, the Russian online social network Odnoklassniki helps people to find their old classmates. All people are also the members of the numerous real social networks like network of friends, business contacts, colleagues, lovers of board games etc.

These property of gathering around similar interest can help to study such particular problem as hidden populations. Examples of hidden populations are drug users, sex workers, jazz musicians. The participants of the hidden population are hard to reach, it means that there is no easy and quick access to the members of these populations. But once the researcher knows some of them he can use this fact that probably they are in contact with another representatives of this population. The individuals of hidden populations from the network around their interest. In this way the subset of this population can be found by "crawling" this graph.

The methods of sampling that use the contacts of individuals in order to find another members of the population are called *chain-referral methods*.

This way of sampling is different from uniform independent sampling. The homophily in the network plus the this fashion of sampling together increase the bias and variance of the estimators.

The structure of the report is following. In the chapter 1 we will introduce the current method of sampling hidden populations that is called respondent-driven sampling (RDS). We will discuss what are the problems and challenges related to this method. Next we will propose the enhancement of the current RDS method that we will study later. The chapter 2 will develop the mathematical model for

studying RDS. Particularly we will introduce the model for controlling the level of homophily for the given network structure. Further using the built model we will find the exact expressions for variance and bias for existing RDS model and enhanced in order to compare them later. In the chapter 3, we will compare the performance of the currently used RDS method to our method. Also we will give some practical recommendations for the users of RDS.

Chapter 1

Respondent-driven sampling

1.1 Motivation

In order to make the statistical analysis about a particular population researchers need to find the representative subset of this population. After analyzing collected data they can produce the results generalized on the whole population.

However sometimes they can have difficulties in finding the subset of the some population. Examples can be the population of the drug users, or the men who have sex with men, or the sex workers, or illegal immigrants, or participants in some social movements, or homeless and so on [18]. The members of these populations are hard to find, they are not willing to participate in the research, that's why these populations are called *hidden populations*.

The reasons to study them can be various. A lot of research is targeted on the studying the HIV prevalence among hidden populations like drug users, female sex workers [15], gay men [17]. Studying the prevalence of HIV can help to understand and control the spreading of the disease. The study [13] was particularly exploring how the connections or the level of professional activity influence more on the income level.

In order to collect samples from the hidden population some special technique were developed. The main challenge here is that the technique should assure that samples are collected with the same probability (or with the known probability). For example, using telephone survey in order to collect information would automatically exclude some subsets of people (like homeless, poor) that don't have telephone number. This fact can affect the correctness of the estimate because it is impossible to predict the bias.

Time-space sampling tries to solve this problem by sampling from different venue-time segments. Each venue-time segment is selected with probability based on the expected number of participants at this segment. However it is impossible to sample exhaustive list of all venue-time segments that results again in the unknown bias of obtained estimation.

Another approach to sample hidden population is using chain-referral techniques. The researchers benefit from the fact that people tend to know each other when they have common interests. A jazz musician has a lot of reasons to know another jazz musicians: they may perform in the same clubs, they probably attend the same events, they may collaborate with each other and so on. The researcher can find just several representatives of hidden population and then ask them to provide the contacts of another members of this population.

One particular chain-referral technique that is currently widely used for studying prevalence

of HIV/AIDS among injection drug users, sex workers, men who have sex with men is called *respondent-driven sampling* (RDS).

1.2 Technique of respondent-driven sampling

RDS begins with selecting group of initial participants that are called *seeds*. The procedure follows according to chain-referral model: each participant in study recruits another participants. The step is called wave. Both participating in the research and recruiting new participants are encouraged by financial incentive. The sampling continues in this way until needed size of participants is reached. During RDS participants are asked to report how many contacts they have. This process enables to collect data for making statistical analysis.

In order to study formally RDS we will model the network of people as the graphs, where the nodes are presented by the individuals and the contact between two individuals is represented by the edge between the correspondent nodes. Then RDS can be regarded as random walk on the graph when the following assumptions are true:

1. Seeds are chosen proportionally to their degree in the network.
2. If individual A knows individual B than individual B knows A as well (network can be represented as undirected graph).
3. The same individual can be recruited multiple times (sampling with replacement).
4. The choice of contacts to recruit is uniformly at random.
5. Individuals know and report precisely their network degree.
6. Each individual is reachable from each other individual (network is connected).

For this process stationary distribution is exactly distribution proportional to network degree. So first assumption guaranties that not only first but all samples during the process are taken with probability proportional to the degree of participants in the network. In [18] this assumption is considered to be reasonable as the people that are drawn as seeds are well-known people and they have usually more contacts than on average. Without this assumption first there should be performed enough number of waves until sample can be considered drawn from stationary distribution.

Some of these assumptions are not so restrictive, like the second one. The other assumptions are arguable. The sensitivity of violation these assumptions were studied in [9] with the simulation experiments.

1.3 Current estimators

Some unbiased estimators were developed in order to generalize the results on all population from the collected samples. We will look at some of them: Sample average (SA), Volz-Heckathorn estimator (VHE) that was introduced in the paper [12] [check this] and estimator presented in the paper [18] that we will call Group estimator (GE).

1.3.1 Sample average

Let y_1, y_2, \dots, y_n be all collected samples during RDS. The simplest estimator of the population mean is just a samples average:

$$\hat{\mu}_{SA} = \frac{y_1 + y_2 + \dots + y_n}{n}$$

But this estimator is biased. In the way that sampling is performed the individuals with more contacts are more likely to be recruited. The probability to encounter the node on the step i with the value y_i is proportional to its degree d_i . So this estimator is biased towards the nodes with higher degrees.

1.3.2 Volz-Heckathorn estimator

In the way that sampling is performed the individuals with more contacts are more likely to be recruited. To correct this bias the responses from individuals are weighted according to their number of contacts. Volz-Heckathorn estimator corrects the bias toward the nodes with higher degrees in the following way.

Let y_1, y_2, \dots, y_n be all collected samples during RDS. Let denote as d_1, d_2, \dots, d_n the number of contacts accordingly to the observed samples. Then estimate μ of the population mean is defined as:

$$\hat{\mu} = \frac{1}{\sum_{i=1}^n 1/d_i} \sum_{i=1}^n \frac{y_i}{d_i}$$

It is shown that this estimator produces asymptotically unbiased results. [put the reference]

1.3.3 Group estimator

This estimator is targeted on the estimation of the percentage of population with certain characteristic. Let's say that the people from the group A possess this characteristic and people from the group B no.

Again each observed individual should report his number of contacts: d_1, d_2, \dots, d_n and whether he/she possess the asked trait. After collecting all the samples the total number of the representatives of the group A , denoted as n_A , and the total number of the representatives of the group B , denoted as n_B , are counted.

Then the estimators of the average number of contacts in the group A and in the group B are respectively:

$$\hat{D}_A = \frac{n_A}{\sum_{i=1}^{n_A} \frac{1}{d_i}}$$
$$\hat{D}_B = \frac{n_B}{\sum_{i=1}^{n_B} \frac{1}{d_i}}$$

Observing the chain of recruitment the number of recruitment between and inside groups are counted. Let's denote r_{AB} as number of rectuitments from the individual in the group A to the individual in the group B , r_{AA} as number of rectuitments from the individual in the group A to the individual in the group A and in the same way r_{BB} , r_{BA} .

Then the estimates for the probability to hire person from the group B by the person from the group A and in the opposite way are respectively:

$$\hat{C}_{A,B} = \frac{r_{AB}}{r_{AA} + r_{AB}}$$

$$\hat{C}_{B,A} = \frac{r_{BA}}{r_{BB} + r_{BA}}$$

Finally the estimate of the fraction of population from the group A is:

$$\widehat{PP}_A = \frac{\hat{D}_B \cdot \hat{C}_{B,A}}{\hat{D}_A \cdot \hat{C}_{A,B} + \hat{D}_B \cdot \hat{C}_{B,A}}$$

The prove that the Group estimator produces asymptotically unbiased results can be access in [18].

Let's note with the VHE it is also possible to estimate the percentage of the population with the given trait. If the i th observed sample obtain this trait then the value if y_i should be set at 1, otherwise at 0.

[to put may be that this is exactly VHE estimator if each person recruits exactly one person]

1.4 Problems of RDS

To understand what are the problems with respondent-driven sampling let's compare it with the "ideal" sampling: independent uniform sampling.

First, let's look at the variance of the uniform sampling. Let y_1, y_2, \dots, y_n be the samples that are taken uniformly at random and all the samples are independent. Let's take the average value of the samples as an estimator of mean population value:

$$\hat{\mu}_{SA} = \frac{y_1 + y_2 + \dots + y_n}{n}$$

This estimator is unbiased and its variance depends on the sample size n . Let σ be the variance of the samples, then the variance of the estimator $\hat{\mu}_{SA}$:

$$\begin{aligned} \sigma_{\hat{\mu}_{SA}}^2 &= \text{var} \left(\frac{y_1 + y_2 + \dots + y_n}{n} \right) = \frac{1}{n^2} \text{var}(y_1 + y_2 + \dots + y_n) = \\ &= \frac{1}{n^2} (\text{var}(y_1) + \text{var}(y_2) + \dots + \text{var}(y_n)) = \frac{1}{n^2} (\sigma^2 + \sigma^2 + \dots + \sigma^2) = \frac{\sigma^2}{n} \end{aligned}$$

In this case the only thing that we can do in order to have better estimation is to increase sample size. To have more precise result we should try to take as much samples as possible. If each sample has a cost, then sample size is restricted by the budget of the research project.

1.4.1 RDS is not uniform

Of course the way we perform sampling it is not independent uniform sampling. First, nodes are not sampled not uniformly. When each participant select another participants with the same probability among his friends there is bias towards the nodes with high degrees. So the more contacts an

individual has more probable he will be invited to participate in the study. In some cases, when we can control the probability of selecting the next participant we can achieve uniform sampling even with random walk.

The study [10] was using Metropolis-Hasting Random Walk to sample Facebook. This method requires information about user's degree and the degrees of all his neighbors. According to this information the selection probabilities are counted for all the friends and then one of them is selected by the computer. In Facebook it is feasible to do: with API requests needed information is collected and then probabilities to choose one of the user's friends are counted. Though this method gives us mechanism to sample uniformly even with random walk, it is not really possible in the situations where we can't control selection probabilities. Like in the case with hidden populations: it is individual who decides how to hire.

Another way to remove degree bias is by reweighing samples according to their degrees, what actually the VHE estimator does, we discussed this in the previous section.

1.4.2 RDS is not independent

Second, collected values are not independent. The participants i and $i + 1$ know each other, they can be friends, relatives or just acquaintances, so their values y_i and y_{i+1} can be dependent. The drug users may be in contact because they go to then same drug dealer and probably buy the same kind of drugs.

The tendency of people with connections to have the similar characteristics is called *homophily*.

And indeed we encounter often a homophily in the real situations. A lot of real networks demonstrate that the value on the node depends from the values of its neighboring nodes. For instance, the study [8] is evaluating the influence of social connections (friends, relatives, siblings) on obesity of people. Interestingly, if a person has a friend who became obese during some fixed interval of time, the chances that this person can become obese are increased by 57%.

Another study [19] that analyzes the data of users and their interactions in the MSN Messenger network found strong relation between users communication behavior (the number of messages exchanged, the total time of chatting, etc) and attributes such as age, gender and even query requests!

On the figure 1.1 the links present the similar votes of U. S. Senators during 2007. With the red labels representing Republicans, the blue labels representing Democrats, the brown labels representing two Independents we can vividly observe the homophily in this network.

Therefore when the values y_1, y_2, \dots, y_n are collected with RDS, then the variance of estimator:

$$\begin{aligned}\sigma_{\hat{\mu}_{SA}}^2 &= \text{var} \left(\frac{y_1 + y_2 + \dots + y_n}{n} \right) = \frac{1}{n^2} \text{var}(y_1 + y_2 + \dots + y_n) = \\ &= \frac{1}{n^2} \left(\sum_{i=1}^n \text{var}(y_i) + 2 \sum_{i=1}^n \sum_{j=i+1}^n \text{cov}(y_i, y_j) \right) = \frac{\sigma^2}{n} + \frac{2 \sum_{i=1}^n \sum_{j=i+1}^n \text{cov}(y_i, y_j)}{n^2} = \\ &= \frac{\sigma^2}{n} \left(1 + \frac{2 \sum_{i=1}^n \sum_{j=i+1}^n \text{cov}(y_i, y_j)}{n\sigma^2} \right) = \frac{\sigma^2}{n} \left(1 + \frac{2 \sum_{i=1}^n \sum_{j=i+1}^n \text{corr}(y_i, y_j)}{n} \right)\end{aligned}$$

So variance of estimator is influenced by some correlation factor $f(n)$ where $f(n) > 1$, that depends on how much the values are correlated:

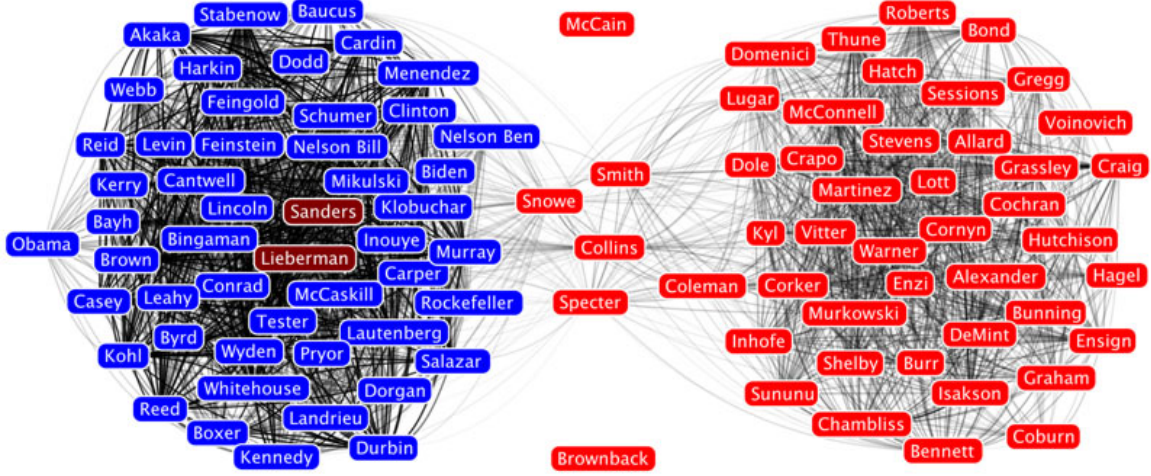


Figure 1.1: Voting patterns of U. S. Senators during 2007 [1]. The red labels represent Republicans, the blue labels represent Democrats, the brown labels represent two Independents.

$$\sigma_{\hat{\mu}_{SA}}^2 = \frac{\sigma^2}{n} f(n)$$

On one hand, when we increase number of participants the factor σ^2/n decreases, but there is also correlation factor $f(n)$, the bigger number of participants, the bigger is this correlation factor. And then in order to improve estimation we could try reduce this correlation factor.

In some literature as in [12] the ratio of the variance of RDS estimate to the variance of estimate obtained from independent uniform random sampling is called design effect, d . They warn the users of RDS: to have the same variance using independent uniform random sampling we need n samples while using RDS need dn samples. We can notice that the correlation factor $f(n)$ is exactly the design effect.

1.4.3 Other problems

RDS can perform poorly if the groups of individuals form different communities. In [11] it is shown that 'bottlenecks' between different groups in hidden population increases variance of RDS estimator. They try RDS on network structure with communities, but where individuals, that are in contact with each other, do not have similar traits and showed that such structure indeed affects on RDS estimate.

1.5 Enhanced RDS

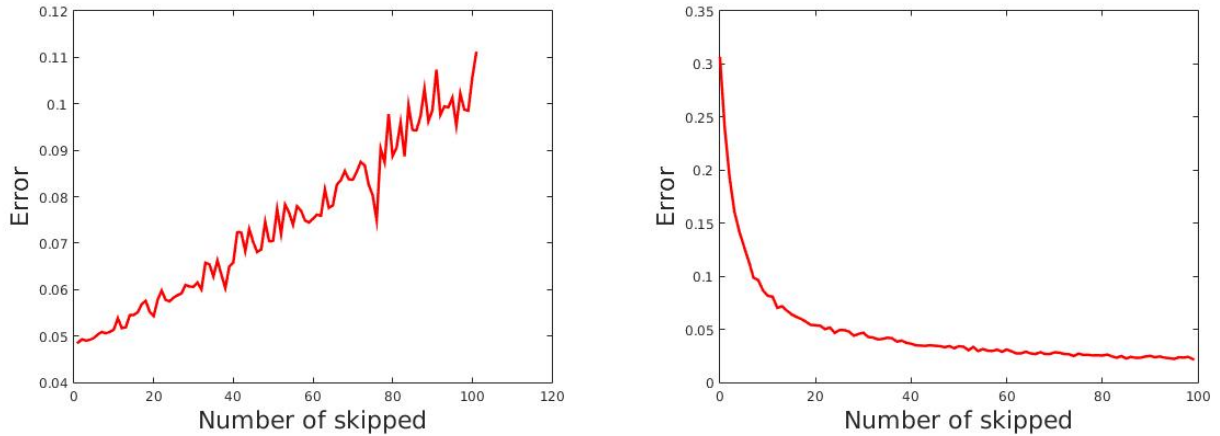
We will state two observations. Combining them we will try to improve current RDS technique.

Observation 1 *Just thinning of sample doesn't help*

In order to reduce correlation between sampled values one could try to thin out sample. It means that instead of taking all collected values y_1, y_2, \dots, y_n into the estimation we can take only, let's say, each second value, y_1, y_3, \dots, y_n . The samples y_1, y_3, \dots, y_n are less correlated then all samples, thus we will reduce correlation factor $f(n)$, so we can expect some improvement. In the same by taking each second values we will degrees in two times the total number of samples. It is not clear if the correlation was reduced enough to compensate this fact. What we observed experimentally (latter we will show it formally) that in general just discarding some samples with the fixed size of recruitment chain will not improve estimation. Experimental results are presented on the figure 1.2 (a).

Observation 2 *Skipping reduces variance*

However if the size of the sample is fixed the further are individuals in the chain of recruitment from each other the better. Let's say that we want exactly 30 participants for study and we have options. We can perform RDS until we reach 30th person and then take each one of them, y_1, y_2, \dots, y_{30} in the estimation. Or we can continue RDS until we reach 60th person and then take each second of them y_2, y_4, \dots, y_{60} in the estimation. The both variants have the same sample size but in the second scenario the values are less correlated. So we expect that the correlation term will decrease. Again on the figure 1.2(b) we can see experimental results and further we will show it formally.



(a) 1000 samples are collected with RDS. The sample size is changed as we skip some values. We can observe that error grows when we try to discard each second, each third and so on nodes.

(b) In all cases the sample size is the same. The number of nodes collected with RDS is different. We observe that error decreases when the distance between samples increases

Figure 1.2: Experiments on the data from Project 90

Keeping in mind that we have fixed budget, the second scenario implies that we don't pay to the participants 1, 3, ..., 59 or we pay to all 60 individuals half of what we payed in the first scenario. But as the motivation to participate in the study is money for the same job people should get the same amount of money.

What we propose lies in the middle. We keep the idea of skipping some values, that will decrease the correlation, but we also suggest to pay less to the individuals that we don't take into the estimation. So among people that are willing to participate some of them will be asked both: to make tests (like blood test, questionnaire) and recruit another participants, let's call them *participants* and

some of them will be asked only to recruit other participants, let's call them *informators*. As the informators make less efforts they will be payed less.

Let's say that each informator receives C_1 units of money and each participant receives $C_1 + C_2$ units of money. The amount of money that we can spend on the research, the budget, is fixed and is denoted as B . Let n be the length of the chain of all the recruitments (it means that informators and participants in total are n).

Let parameter k be the number of skipped samples in the chain between participants (it means that we take only each k th individual as a participant participant, the rest $k - 1$ individuals are just informators).

The reason to use informators is to try to reduce correlation by making bigger the distance between participants and therefore less correlation. If the informators were willing to do their job for free then, according to the third observation, we would try to have as more informators between participants as possible. But all the idea of respondent-driven sampling holds on the money incentives, without payment nobody will do anything. For this reason informators should be also payed, but less then actual participants. As we still spend part of the budget on informators the number of participants will decrease with increasing number of informators.

To understand better the idea, let's imagine that we have 60 euros budget, each informator is payed 2 euros, $C_1 = 2\text{€}$, each participant is payed 10 euros, $C_2 + C_1 = 10\text{€}$. On the figure you can see different scenarios of RDS for the same budget.

On the figure 1 the parameter k is equal to 1. It means that everybody is participant and payed 10 euros. Then with our budget we can collect 6 samples that due to homophily can be highly correlated.

On the figure 2 the parameter k is equal to 2. It means that we take one participant, one informer, one participant and so on. We can see clearly that we went deeper in the network, the recruitment chain is longer. In this case on our budget we can collect only 5 samples, but they are less correlated than in the previous situation.

Both scenarios require the same budget, so in the terms of money they are equal. But what is better 6 more correlated samples, or 5 but less correlated samples?

If the characteristic that we try to estimate does not correlate between friends or people who have contact then it is useless to discard some values (we just pay for nothing). But if the values are highly correlating intuitively skipping can help a lot. There is a trade-off: on one hand we make the chain longer and reduce dependency between participants. On other hand we still spend money on people who do not bring any information needed for research and finally there will be less participants.

The better scenario will be the one that has the least error. If we are able to quantify the error depending on the number of samples that we skip $k - 1$, we will be able to suggest the sampling method that will bring the most precise result. That is the question that we are going to analyze formally in the next chapter.

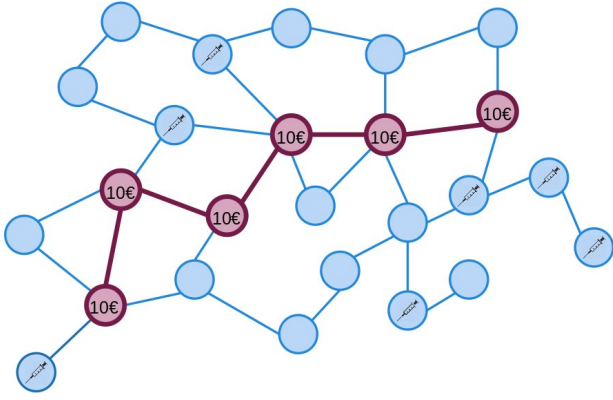


Figure 1.3: Scenario 1

Budget B	60€
C_1	2€
C_2	8€
k	1
Sample size	6
Chain size	6

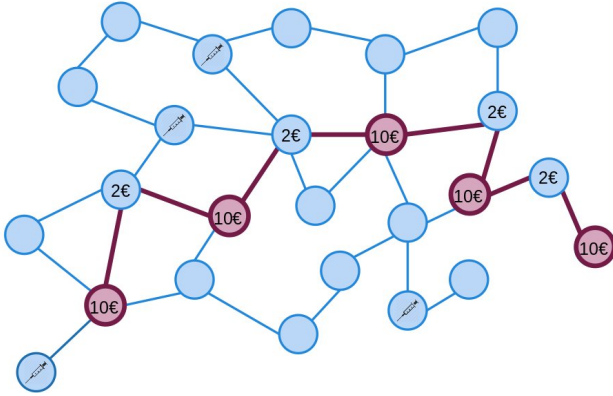


Figure 1.4: Scenario 2

Budget B	60€
C_1	2€
C_2	8€
k	2
Sample size	5
Chain size	9

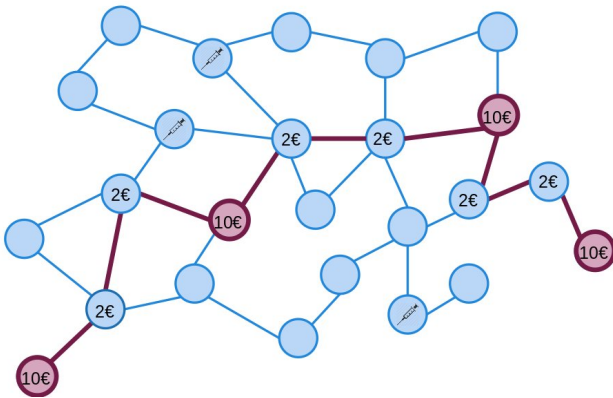


Figure 1.5: Scenario 3

Budget B	60€
C_1	2€
C_2	8€
k	3
Sample size	4
Chain size	10

Chapter 2

Mathematical model

2.1 Network modeling

As it has been already mentioned we will model with the graph the network formed by the members of hidden population. The nodes of this graph correspond to the persons. The edge between two nodes represents the connection between corresponding persons.

We have multiple possibilities to obtain graph structure. First, there are multiple available network structures that were collected from the real networks. Particularly the Stanford Large Network Dataset Collection [6] provide the data of online-social networks (we will use part of Facebook graph), collaboration networks, web graphs, Internet peer-to-peer network and a lot of others.

The other possibility is to use random graph theory to generate the random graph that will imitate the network structure. There is number of different random graphs. Particularly we used Erdos Renui graph, random geometric graph, preferential attachments graph.

2.2 Network with values

2.2.1 Motivation

Each node should maintain the value of the attribute, that is going to be estimated. For instance, if we have a social network the attribute can be the age, gender of a user.

All the random graph models give us possibility to generate only the structure of a network. The next step is to generate the values on the nodes of the obtained graph which will represent needed attribute.

The simplest idea is to assign values randomly to the nodes independently of the graph structure. For example, we could assign the age of the user according to the uniform distribution or normal distribution or any distribution we want our attribute to be distributed. This approach has an explicit weakness: it does not take into account the homophily, the tendency of people with connections to have the similar characteristics.

We already explained that this property is frequently encountered in the real networks.

The reason why we don't want to ignore homophily is because the way of performing sampling and the property of homophily together influence on the sampling variance and bias. Further we will count formally the sampling variance for given network and attribute of the nodes.

So for study purposes we would like to assign values to the nodes of the network in such a way that the value of the node depends on values of its neighbors. The other point is that the level of correlation can be different within the same network but for different attributes.

So moreover we would like to control the level of homophily in the network.

The study [16] was investigating how the binge drinking is influenced by the position of student in the network of students. The students were labeled according to belonging to one the group: member of a binging group, liasons, isoletes, etc. The researchers looked for the relation of the episodes of binge drinking per fixed period of time and the student's label. They found strong dependency while the students were young, but not when they became adults. So for the same network the different attributes: binge drinking in school and binge drinking after school have different level of dependency of the friend's behavior.

Regarding what was said above we would like also to be able to tune the correlation in the network. For this reason the next technique was developed.

2.2.2 Definitions

We have graph with n nodes. To each node i will be assigned random value X_i from the set of values $V, V = \{1, 2, 3, \dots, k\}$.

Instead of looking on distributions of the values on nodes independently, we will look at the joint distribution of values on all the nodes. The distribution should take into account the values of the node's neighbors as well.

Let's denote (X_1, X_2, \dots, X_n) as \bar{X} . We will call \bar{X} as a random field. When random variables X_1, X_2, \dots, X_n take values x_1, x_2, \dots, x_n respectively we will call (x_1, x_2, \dots, x_n) a configuration of a random field and we will denote it as \bar{x} . As the basement of distribution we will take Gibbs distribution, that originally comes from physics [put a reference].

For simplicity, instead of writing $p(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$ where $x_1, x_2, \dots, x_n \in V$ we will write $p(\bar{X} = \bar{x})$ or just $p(\bar{x})$.

For each possible configuration \bar{x} we will associate the number that is called global energy of the graph that is counted in the following way:

$$\varepsilon(\bar{x}) = \sum_{i \sim j, i \leq j} (x_i - x_j)^2$$

where $i \sim j$ means that the nodes i and j are neighbors in the graph.

Let's turn our attention to the one node i . We will define local energy on the node i as:

$$\varepsilon_i(x_i) = \sum_{j|i \sim j} (x_i - x_j)^2$$

Then we can rewrite the expression of the global energy knowing the local energies on all the nodes:

$$\varepsilon(\bar{x}) = \frac{1}{2} \sum_i \varepsilon_i$$

2.2.3 Gibbs distribution

Now let's consider the following probability distribution:

$$p(\bar{x}) = \frac{e^{-\frac{\varepsilon(\bar{x})}{T}}}{\sum_{\bar{x}' \in |V|^n} e^{-\frac{\varepsilon(\bar{x}')}{T}}} \quad (2.1)$$

where T is temperature, $T > 0$.

The reason why it is interesting to look at this distribution follows from the theorem [Theorem 2.1, p. 260], [7]. When random field has distribution 2.1 then the probability that the node has particular value depends only on the values of its neighboring nodes and does not depend on the values of all other nodes.

Let's denote N_i as the set of neighbors of the node i . If the L is the subset of nodes then X_L will denote the set of random variables of the corresponding nodes. The last property can be formulated in the following way:

$$p(X_i = x_i | X_{N_i} = x_{N_i}) = p(X_i = x_i | X_{\{1,2,\dots,n\} \setminus i} = x_{\{1,2,\dots,n\} \setminus i})$$

This property is called Markov property.

Moreover for each node i , knowing values of its neighbors, we can write the distribution of values: as following:

$$p(X_i = x_i) = \frac{e^{-\frac{\varepsilon_i(x_i)}{T}}}{\sum_{x' \in V} e^{-\frac{\varepsilon_i(x')}{T}}}$$

The temperature parameter T plays very important role of the tuner of the correlation level in the network. Later we will show some examples for better understanding.

Gibbs distribution found many interesting applications in real-world problems. Particularly it lies in the basement of the proposed in [14] distributed algorithm for channel selection of the Access Points. The channels should be selected in such way that interference in the network is minimized.

2.2.4 Algorithm

Practically speaking direct sampling from the distribution 2.1 is not so easy. Let's just notice that the number of possible configuration \bar{x} is $|V|^n$, where $|V|$ is size of the values set and n is number of the nodes in graph, as to each from n nodes we need to assign the value from the set V . In this way for the graph with just 100 nodes and 10 possible values it would make up 10^{100} possible configurations. Then probability for each of 10^{100} possible configurations should be counted, which would require huge precision from the computer to be able to sample from derived distribution.

In order to produce samples from this distribution we are using Gibbs sampler [7].

Let's regard each configuration \bar{x} as a state of the Markov chain. Let's denote as \bar{x}^k the state at the step k . There is positive probability to transit from one state to another if the corresponding configurations differ only in the value of one node and the values of all other nodes are the same. Let's take two configurations $(x_1, x_2, \dots, x_i, \dots, x_n)$ and $(x_1, x_2, \dots, x'_i, \dots, x_n)$ which differ only in the value of the node i . Whatever the value was on the node i in the first configuration the probability to transfer from the first state to the second is:

$$\frac{e^{-\frac{\varepsilon_i(x'_i)}{T}}}{\sum_{x' \in V} e^{-\frac{\varepsilon_i(x')}{T}}}$$

Interestingly, the stationary distribution of this Markov Chain is exactly 2.1.

Therefore following algorithm after converging will produce a sample from the distribution 2.1.

1. Create random configuration of properties on all nodes.
2. Choose the node i
 - according to some distribution $q = q_1, \dots, q_n$ or
 - visiting each node consequently (periodic Gibbs sampler)
3. For each value $x \in V$ count the local energy on chosen node i as

$$\varepsilon_i(x) = \sum_{j|i \sim j} (x - x_j)^2$$

4. Choose a new value x_i according to probability

$$\frac{e^{-\frac{\varepsilon_i(x_i)}{T}}}{\sum_{x' \in V} e^{-\frac{\varepsilon_i(x')}{T}}} \quad (2.2)$$

where T is temperature.

5. Continue 2-3 needed number of iterations.

2.2.5 Explanatory example

To understand the influence of the temperature on the value distribution 2.2 we will look at the following example.

Let's say that we have the graph to which nodes we want to assign values 1, 2, 3, 4, 5. Now let's look only at the vertex A its neighbors B, C, D, E which have assigned values 1, 5, 3, 4 respectively (figure 2.1). And now it is turn of A to choose a value.

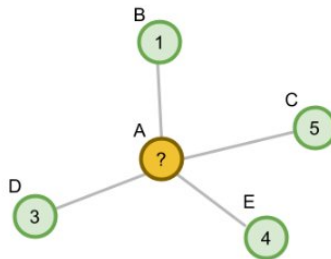


Figure 2.1: Node A and his neighbors

With different values the node A will have different local energy. Let's summarize it in the table.

Value on the node A	1	2	3	4	5
Corresponding energy	29	15	9	11	21

We can see that with different values the local energy on the node A will be different. And according to the distribution 2.2 the values that bring low energy is more preferable: the less is the energy that causes particular value the more probable it will be chosen. In this example the lowest temperature corresponds to the values 3, so it will the highest chance to be picked.

It is not only the energy that influences on the distribution 2.2. There is also temperature parameter T . To feel the impact of temperature we will present the distribution of values for different temperature T in the next table.

Temperature	$p(A = 1)$	$p(A = 2)$	$p(A = 3)$	$p(A = 4)$	$p(A = 5)$
0.1	0.0000	0.0000	1.0000	0.0000	0.0000
1	0.0000	0.0022	0.8789	0.1189	0.0000
10	0.0483	0.1957	0.3566	0.2920	0.1074
100	0.1769	0.2035	0.2161	0.2118	0.1917
10000	0.1998	0.2000	0.2002	0.2001	0.1999

We can see that when the temperature is 0.1 the probability to choose the value 3, that has the lowest energy, is 1. So the values of the neighbors B, C, D, E indeed impact a lot on the value of the node A . In this case, the value of A is dictated by the values of its neighbors. This is exactly what we wanted to achieve: that characteristic of the person is influenced by its contacts. However we may want to have such dependency but not so high.

We can try to play with the temperature parameter. As the temperature increases we can observe more "randomness". Thus when the temperature is 1 the value 3 is still very probable, but there is also some positive probability that values 2, 4 will be picked. And the more we increase temperature the higher becomes probability to choose the value different from 3. We can interpret this as the person still depends on the connections but he has some "free choice".

When the temperature is really high the choice of value will not almost depend on the values of its neighbors.

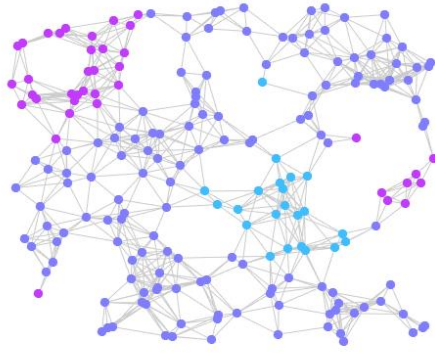
With this example we observed that the distribution favors the values that bring local energy of the node to minimal and the lower the temperature is the more favorable they are. So the higher is temperature the less is the decency between value of the person from his contacts.

2.2.6 Demonstration of random graphs with values

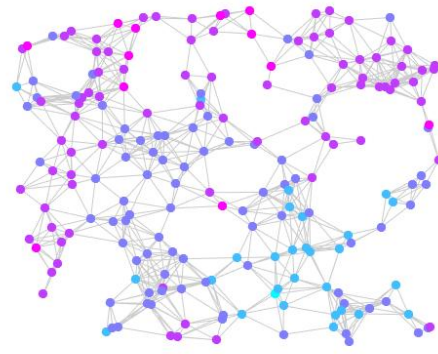
In order to demonstrate that proposed model works first we will show the generated graphs with values to see the result visually and then we will look at the ways to measure level of values dependency in the graph.

On the figure 2.2 presented the same random geometric graph with 200 nodes and radius 0.13, $RGG(200, 0.13)$ where the values $V = \{1, 2, \dots, 5\}$ are chosen according to the Gibbs distribution. The values are depicted on the pictures as colors.

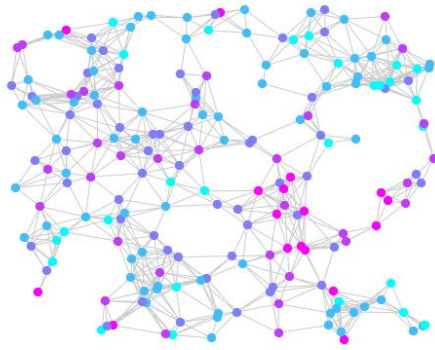
From the pictures we can observe that the level of dependency between values of the node changes with different temperature. When temperature is 1 we can distinctly distinguish clusters. With increasing temperature, 5 and 20, the values of neighbors are still similar but with more and more variability. When temperature is very high then the values seem to be assigned randomly.



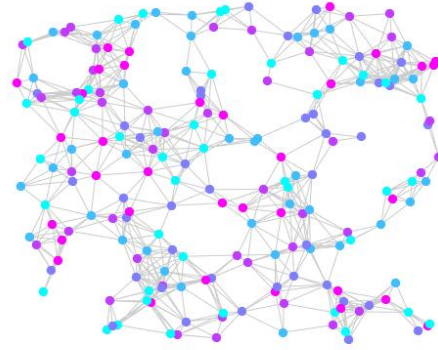
(a) Temperature 1



(b) Temperature 5



(c) Temperature 20



(d) Temperature 1000

Figure 2.2: RGG(200, 0.13) with generated values for different temperature

the correlation of the gibbs field??? maybe to put but then say that it is not the same

In order to give more formal illustration we can look at the correlation between values of the nodes that we see during the random walk on the graph.

Let's denote as $Y_0, Y_1, \dots, Y_i, \dots$ the sequence of values on the nodes that we observe during the random walk. Let's say that we start random walk with stationary distribution, so the values $Y_0, Y_1, \dots, Y_i, \dots$ have the same stationary distribution. Then covariance between two values Y_i and Y_{i+k} depends only on the distance k in the sequence $Y_0, Y_1, \dots, Y_i, \dots, Y_{i+k}, \dots$ between them:

$$\text{cov}(Y_i, Y_{i+k}) = \text{cov}(Y_0, Y_k)$$

explain why

On the figure 2.3 we present correlation between values depending on this distance k for the graphs shown above. In this way we can assure us that the higher is temperature the less correlated values of neighbors are.

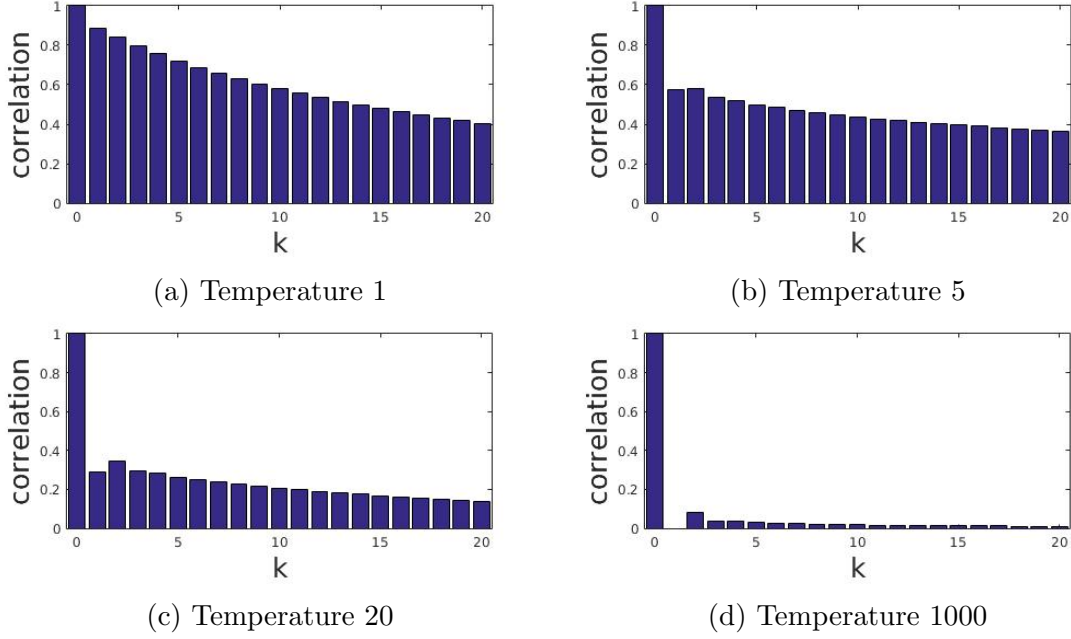


Figure 2.3: Correlation of the values of the nodes observed during the random walk depending on the difference in their order for different temperature

2.3 Expected energy in steady state

2.3.1 Motivation

The question that is still not clear about the algorithm is when to stop it. How many steps are enough to perform in order to claim that achieved configuration is indeed sampled from the Gibbs distribution? In order to answer this question we should try to understand if we can detect somehow that the process described in the subsection 2.2.4 has reached its stationary state.

For this purpose we can try to use the notion of the global energy that can be counted for each configuration. Maybe it can signalize us when it is save to stop the algorithm.

Let's look how the global energy of the graph is changing during the steps of algorithm. For each configuration \bar{x} the global energy of the graph is:

$$\varepsilon(\bar{x}) = \sum_{i \sim j, i \leq j} (x_i - x_j)^2$$

For each iteration of the algorithm when the configuration is changed the global energy is counted. The results are presented on the figure 2.4 .

At the beginning of the algorithm values are assigned to all nodes uniformly from all possible values. We can observe on the figure 2.4 that the energy for this first random configuration is the highest. As values on the nodes are changing according to Gibbs sampling the energy decreases (with some variation). After about 200 steps (it means that up to this time each node updated its state once) we can see that the changes of energy do not exceed some thresholds. So after some time energy comes to some value, stabilizes and does not change a lot. We will call this value as **expected global energy**. Knowledge about this value and its variation can indicate us when it is time to finish the algorithm.

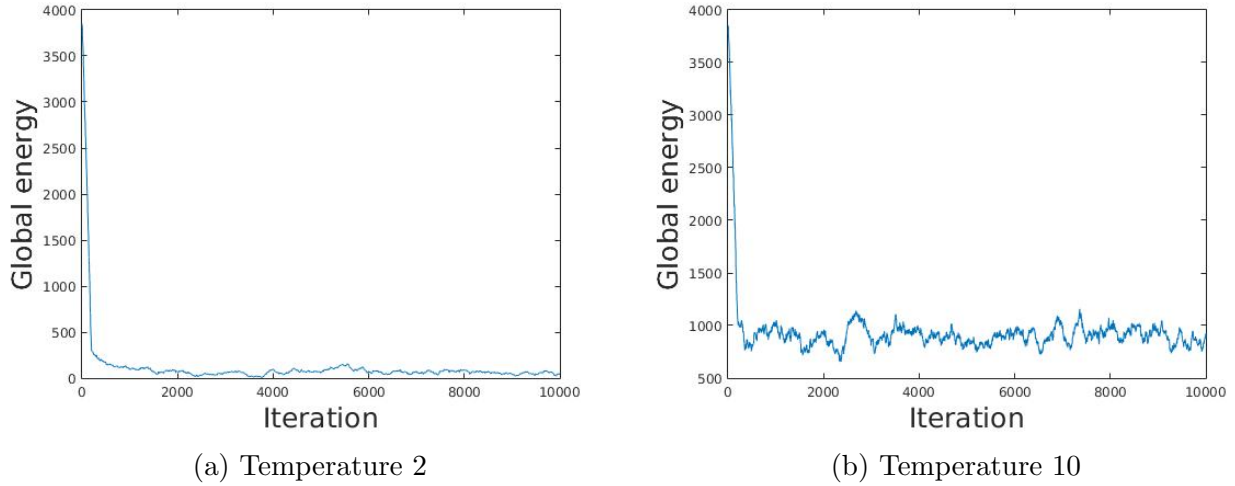


Figure 2.4: Energy changing with iterations of the algorithm

There is another reason why we would like to predict energy. We saw previously that by varying the temperature T we can change how strongly values of the neighbors are related: low temperature brings high correlation of values and high temperature brings almost random assignment of values. However it is impossible to use only temperature T as a metric of values correlation. Temperature by itself does contain a lot of information. For the graph $\text{RGG}(200, 0.13)$ the temperature 200 we can consider as high (because values are not really correlated) but for the graph $\text{RGG}(2000, 0.06)$ it is not the case. We can't judge the level of correlation only by temperature, we should take into account number of nodes, number of edges, structure of graph, possible values that can be assigned to nodes as well.

That's why we need another metric. And again the knowledge about expected global energy can help us. On the figure 2.4 we can see two illustrations of energy changing with time for different temperature. Moreover for different temperature there is different expected global energy (and the its variance is also different). On the both pictures starting from the same initial configuration but applying the algorithm with different temperatures we come to different expected energy.

Then the more appropriate metric can be following number: in how much times energy decreases from its initial value (that corresponds to the random configuration) to its stable value (the one that we want to be able to predict).

The problem here is that the algorithm needs the temperature parameter. So after deciding that we want to decrease the initial energy in 5 times we still need to understand which temperature will bring the system to this target energy.

For this purpose we were interested in finding dependency of the global energy from the temperature.

2.3.2 Analysis

For the reasons discussed above we would like to know the expected value of the global energy.

Let the random variable $\varepsilon(\bar{X})$ be total energy of the graph with random field $\bar{X} = (X_1, X_2, \dots, X_n)$. We can write expected energy by definition as:

$$E[\varepsilon(\bar{X})] = \sum_{\bar{x}} p(\bar{x}) \cdot \varepsilon(\bar{x}) \quad (2.3)$$

where probability of one particular configuration \bar{x} is

$$p(\bar{x}) = \frac{e^{-\frac{\varepsilon(\bar{x})}{T}}}{\sum_{\bar{x}' \in |V|^n} e^{-\frac{\varepsilon(\bar{x}')}{T}}} \quad (2.4)$$

Both in 2.3 and in 2.4 formulas summation is over all possible configurations and the number of all possible configurations is huge, $|V|^n$. One of the ways to calculate such kind of expressions would be following: using Gibbs sampling, run the algorithm until convergence large enough amount of times, and then average the result. But actually it is the opposite of what we are trying to do. In this way we can build empirical dependency of expected energy from the temperature. But to have these empirical results we need to perform the algorithm large amount of times and for different temperatures. Such simulations can take a lot of time, especially because we don't know when it is safe to stop algorithm. So actually we would like to have some theoretical results (at least approximated, just to have a notion about energy dependency from the temperature).

First, let's rewrite global energy as:

$$\varepsilon(\bar{X}) = \sum_{i \sim j, i \leq j} (X_i - X_j)^2$$

Then the expected global energy of a field is:

$$E[\varepsilon(\bar{X})] = \sum_{i \sim j, i \leq j} E[X_i - X_j]^2$$

In order to count this expression we need to know the distribution of values for each node, X_1, X_2, \dots, X_n . Moreover, we need to know the correlation between random variables X_i, X_j for all pairs i, j .

In reality both these requirements are difficult to satisfy.

In fact, it seems that we know the distribution of values on the nodes. So the first demand should be easy. We have already written that the values on the node i are distributed in the following way:

$$p(X_i = x) = \frac{e^{-\frac{\varepsilon_i(x)}{T}}}{\sum_{x' \in V} e^{-\frac{\varepsilon_i(x')}{T}}} \quad (2.5)$$

But if we look closer at the local energy on the node i , $\varepsilon_i(x)$, we will notice that it implies that the values on the neighboring nodes are known:

$$\varepsilon_i(x) = \sum_{j|i \sim j} (x - x_j)^2$$

It means that distribution 2.5 is actually conditional on the values of the neighbors. We can

rewrite it as:

$$p(X_i = x) = \frac{e^{-\frac{\varepsilon_i(x)}{T}}}{\sum_{x' \in V} e^{-\frac{\varepsilon_i(x')}{T}}} = \sum_{a_1, \dots, a_j \subset V} \prod_{j|j \sim i} p(X_j = a_j) \frac{e^{-\frac{\sum_{j|j \sim i} (x - a_j)^2}{T}}}{\sum_{x' \in V} e^{-\frac{\sum_{j|j \sim i} (x' - a_j)^2}{T}}}$$

We can see that this expression includes also probabilities for other nodes to have some particular values. So the probability to have a value on the node i depends on the values of its neighbors that are also dependent from their neighbors and so on.

The second demand of knowing correlation for all pairs X_i, X_j require to know the joint distribution of these random variables. That brings us to the same problem: no way to write explicitly the expression of this distribution

That's why in order to calculate expected energy at least approximately we will make some assumptions.

The same expected value

The experiments showed us that the expected value of the each variable X_1, X_2, \dots, X_n is the same and equals to the average value of the values from the set V .

$$E[X_i] = av_V$$

We didn't show it formally, that's why we will write that this is an assumption. But we have reasons to believe that it is true and can be proved.

Assumptions about no correlation between neighbors

We will make an assumption that values are assigned independently of each other. Of course it is not true. The distribution of the values on a node takes into the account the values of his neighbors. By making this assumption we will make some error, but we can get the approximated expression. If the approximation is close to the reality it can give us at least the idea about energy-temperature dependency.

Let N_i be the number of the neighbors of the node i . So if there is no correlation between value assigned to the nodes, then we can write global energy as:

$$\begin{aligned} E[\varepsilon(\bar{X})] &= \sum_{i \sim j, i \leq j} E[X_i - X_j]^2 = \frac{1}{2} \sum_{i \sim j} E[X_i - X_j]^2 = \\ &= \frac{1}{2} \sum_{i \sim j} (Var(X_i - X_j) + E[(X_i - X_j)]^2) = \frac{1}{2} \sum_{i \sim j} (Var(X_i - X_j) + (E[X_i] - E[X_j])^2) = \\ &= \frac{1}{2} \sum_{i \sim j} (Var(X_i) + Var(X_j) - 2Cov(X_i, X_j)) = \frac{1}{2} \sum_{i \sim j} (Var(X_i) + Var(X_j)) = \\ &= \frac{1}{2} \sum_{i \sim j} 2N_i Var(X_i) = \sum_{i \sim j} N_i Var(X_i) \end{aligned} \tag{2.6}$$

Special case

Let's look at the case when the values are assigned to the nodes according to the same distribution in independent way (that means that there is no correlation between assigned values). Let m be the number of edges in the graph. Then the expression for energy becomes:

$$E[\varepsilon(\bar{X})] = \sum_{i \sim j} N_i \text{Var}(X_i) = 2m \text{Var}(X_i)$$

where X_i and X_j are random variables with the same distribution. Interesting fact that it reminds us the famous formula for energy!

Particularly, we can compute expected energy in this way for initial random configuration where the values are assigned to the nodes independently and uniformly from all possible values in V . When X_i is distributed uniformly in V the variance of X_i is:

$$\text{var}(X_i) = \frac{|V|^2 - 1}{12}$$

Then the expected energy of the graph on random configuration \bar{x} is

$$E[\varepsilon(\bar{X})] = m \frac{|V|^2 - 1}{6}$$

Assumptions about values distribution

We saw that it is impossible to write explicitly the distribution of the values on one node explicitly. In order to simplify the expression for values distribution on the node i we will also make some assumptions. First, let's say that all neighbors of the node i have the value av_V , $av_V = \text{average}(V)$. If for all $j \in N_i : X_j = av_V$ then $p(X_j = av_V) = 1$. With this assumption probability that the node i will have value $x \in V$ is

$$p(X_i = x) = \frac{e^{-\frac{N_i(x - av_V)^2}{T}}}{\sum_{x' \in V} e^{-\frac{N_i(x' - av_V)^2}{T}}} \quad (2.7)$$

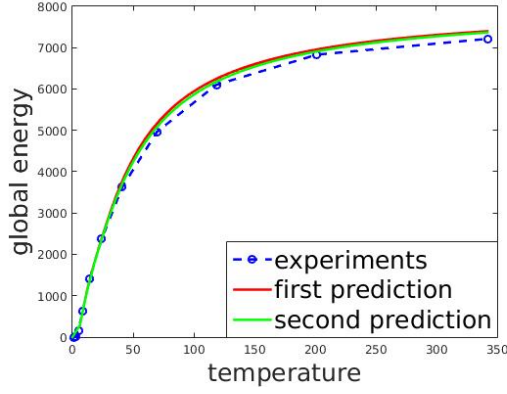
Then combining expressions 2.6 and 2.7 we can write expected energy as:

$$E[\varepsilon(\bar{X})] = \sum_{i \sim j} N_i \text{Var}(X_i) = \sum_{i \sim j} N_i (E[x_i]^2 - (E[x_i])^2) = \sum_{i \sim j} N_i \left(\frac{\sum_{i \in V} i^2 e^{-\frac{N_i(i - av_V)^2}{T}}}{\sum_{x' \in P} e^{-\frac{N_i(x' - av_V)^2}{T}}} - av_V^2 \right)$$

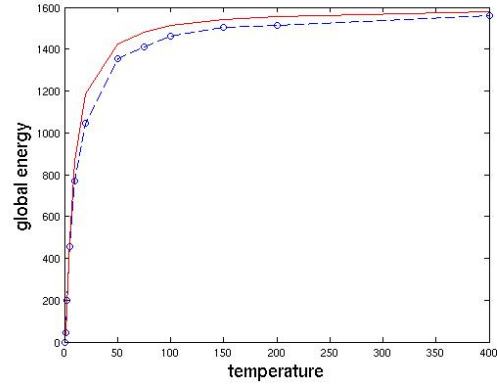
To simplify even more we can assume that each node has the same following distribution of values:

$$p(X_j = x) = \frac{e^{-\frac{d(x - av_P)^2}{T}}}{\sum_{x' \in P} e^{-\frac{d(x' - av_P)^2}{T}}} \quad (2.8)$$

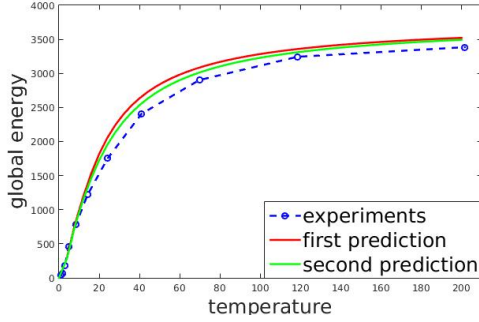
where d is average degree of the graph.



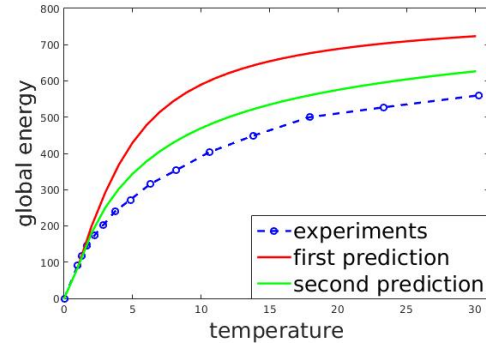
(a) Random ER graph with 200 vertices and values $[1, \dots, 5]$ $p = 0.1$



(b) Grid on torus graph with $200 = 20 \times 10$ vertices and values $[1, \dots, 5]$



(c) Random geometric graph with 200 vertices, radius 0.13 and values $[1, \dots, 5]$



(d) Preferential attachment graph with 200 vertices, 1 link for new arriving node and values $[1, \dots, 5]$

Figure 2.5: Energy prediction for different graphs

Then expected energy is counted in the following way:

$$E[\varepsilon(\bar{x})] = 2mVar(x_i) = 2m \left(\frac{\sum_{i \in P} i^2 e^{-\frac{d(i-av_P)^2}{T}}}{\sum_{x' \in P} e^{-\frac{d(x'-av_P)^2}{T}}} - av_P^2 \right)$$

2.3.3 Results

On the figure 2.5 we can observe predicted energy with assumptions including that the values on the nodes are distributed as in 2.7 (it is called first prediction), predicted energy with assumptions including that the values on the nodes are distributed as in 2.8 (it is called second prediction) and the expected energy counted with simulations.

We can see that both predictions coincides well with the experiments and the first prediction is more accurate.

Even if the result is approximated it can give us great intuition about energy-temperature dependency.

Practically speaking, it takes 20 minutes of experiments to understand that we need to use temperature 30 in order to reduce the initial temperature in 4 times in the presented ER graph just with 200 nodes and just few seconds with obtained result.

2.4 Error prediction

Now, when we have network, where each node maintain the value we can begin to study formally described in the section 1.5 method of enhanced RDS.

Due to homophily in the network and the way the respondent-driven sampling is performed the variance of the estimator will be different from the case if it was just independent uniform sampling. Studying variance is important for multiple reasons. It is essential for building confidence intervals. It is the factor that influences on the error of estimator. When we have multiple estimators and we know the error of each of them we have the instrument to compare them.

In this section we will take the sample average (SA) as the estimator for the population mean. In order to decrease the variance of the estimator we will increase the distance between samples as it was described in the section 1.5. We will look at the different possibilities of sampling hidden populations in the conditions of limited budget. For each scenario we will have different estimators. We are going to choose the best of them, the one that has the minimum error. Then we will compare it with existing estimators.

2.4.1 Variance prediction

Let $Y_1, Y_2, Y_3, \dots, Y_n$ be the samples that are taken during the random walk. Let's take the average value of the samples as an estimator of mean population value:

$$\hat{\mu}_{SA} = \frac{Y_1 + Y_2 + \dots + Y_n}{n}$$

We noted that the variance of the estimator $\hat{\mu}_{avg}$ is influenced by the correlation factor as the random variables $Y_1, Y_2, Y_3, \dots, Y_n$ are correlated:

$$\sigma_{\hat{\mu}_{SA}}^2 = \frac{\sigma^2}{n} f(n)$$

Let's look what will happen with variance in the suggested enhanced RDS.

To remind the notation, we denote budget as B , n individuals that we see during the RDS receive C_1 units of money for recruiting another individuals, each k th person from them take also part in the study and receives C_2 additional units of money.

In this way $\frac{n}{k}$ individuals from n are participants. Then the following equality should be true:

$$B = nC_1 + \frac{n}{k}C_2$$

From this equality we can see that with the fixed budget B and skipping $k - 1$ individuals between participants, the length of the chain n can be:

$$n = \frac{kB}{kC_1 + C_2}$$

When each k th person is a participant, it means that for real we have only values $Y_1, Y_k, Y_{2k}, \dots, Y_{\frac{kB}{kC_1+C_2}}$ that we can take for the estimation. In this way the number of participants m is:

$$m = \frac{B}{kC_1 + C_2}$$

Then for each scenario with the fixed budget B depending on the k the estimator is:

$$\hat{\mu}_{SA}^k = \frac{Y_k + Y_{2k} + \dots + Y_{\frac{kB}{kC_1+C_2}}}{\frac{B}{kC_1+C_2}}$$

We can see that the number of participants depends on k and the bigger is k the less participants we have. Also we observed experimentally that the bigger is k (with the fixed sample size) the less is the correlation. So actually correlation factor depends on k . Saying all this we can write the variance of estimator as:

$$\sigma_{\hat{\mu}_{SA}^k}^2 = \frac{\sigma^2}{\frac{B}{kC_1+C_2}} f(k)$$

Now we have to find the expression for the correlation factor $f(k)$.

2.4.2 Geometric correlation

First we will find correlation for simple example and then we will generalize it for any kind of graph.

For now we will not care about graph structure. Let's assume that our collected samples Y_1, Y_2, \dots, Y_n are correlated in the known way:

$$\text{corr}(Y_i, Y_{i+h}) = \rho^h$$

So correlation between the nodes that are at the distance 1 in the chain have correlation ρ , at distance 2 have correlation ρ^2 and so on. We will refer to this model with indicated correlation as to *geometric model*. Then we can write the variance of the mean estimator as:

$$\begin{aligned} \sigma_{\hat{\mu}_{SA}}^2 &= \text{var} [\bar{Y}] = \text{var} \left[\frac{Y_1 + Y_2 + \dots + Y_n}{n} \right] = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \text{cov}(Y_i, Y_j) = \\ &= \frac{\sigma^2}{n^2} (n + 2(n-1)\rho + 2(n-2)\rho^2 + \dots + 2 \cdot 2\rho^{n-2} + 2 \cdot 1\rho^{n-1}) = \\ &= \frac{\sigma^2}{n^2} \left(n + 2 \sum_{i=1}^{n-1} (n-i)\rho^i \right) = \frac{\sigma^2}{n^2} \left(n + 2n \sum_{i=1}^{n-1} \rho^i - 2 \sum_{i=1}^{n-1} i\rho^i \right) = \\ &= \frac{\sigma^2}{n} \left(n + 2n \frac{\rho - \rho^n}{1 - \rho} - 2\rho \sum_{i=0}^{n-2} (\rho^{i+1})' \right) = \\ &= \frac{\sigma^2}{n} \left(n + 2n \frac{\rho - \rho^n}{1 - \rho} - 2\rho \left(\frac{\rho - \rho^n}{1 - \rho} \right)' \right) = \\ &= \frac{\sigma^2}{n} \left(n + 2n \frac{\rho - \rho^n}{1 - \rho} - 2\rho \frac{(1 - n\rho^{n-1})(1 - \rho) + \rho - \rho^n}{(1 - \rho)^2} \right) = \end{aligned}$$

$$= \frac{\sigma^2}{n^2} \frac{n - n\rho^2 - 2\rho + 2\rho^{n+1}}{(1 - \rho)^2}$$

In the end we have the following expression:

$$\text{var} [\bar{Y}] = \frac{\sigma^2}{n} \frac{1 - \rho^2 - 2\rho/n + 2\rho^{n+1}/n}{(1 - \rho)^2} \quad (2.9)$$

From here we can get that correlation factor is:

$$f(n) = \frac{1 - \rho^2 - 2\rho/n + 2\rho^{n+1}/n}{(1 - \rho)^2}$$

It can be shown that this factor $f(n)$ is increasing function of n , ($n > 1$) and it has its minimum 1 when $n = 1$. It is clear, when there is only one individual there is no correlation, because there is only random variable Y_1 . When new participants are invited, the correlation increase due to homophily as we explained earlier.

Let's look what happens to the correlation factor when n goes to infinity:

$$f(n) = \frac{1 - \rho^2 - 2\rho/n + 2\rho^{n+1}/n}{(1 - \rho)^2} \xrightarrow{n \rightarrow \infty} \frac{1 - \rho^2}{(1 - \rho)^2} = \frac{1 + \rho}{1 - \rho}$$

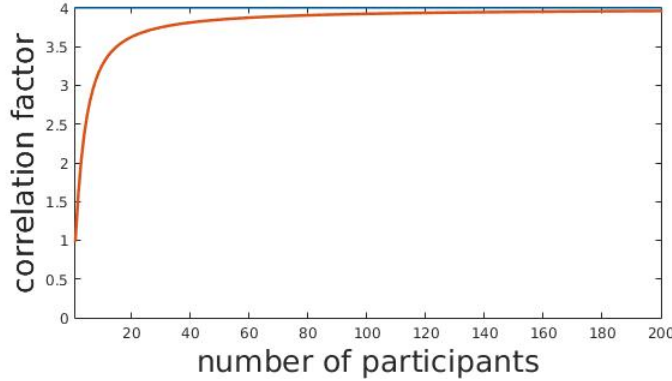


Figure 2.6: Correlation factor depending on the number of the participants when ρ is 0.6

Using the following approximation the expression for the sample variance becomes much more easier:

$$\text{var}_{approx} [\bar{Y}] = \frac{\sigma^2}{n} \frac{1 + \rho}{1 - \rho}$$

Approximation is close when n is big enough. To compare approximated expression with original one, look at the figure 2.7 where parameter ρ is 0.6. As it is reasonable to suppose that the sample size is bigger than 50, we can consider this approximation good enough in this case. The reason to use this approximation is that the expression is much simpler and some facts that are very difficult to prove for real $\text{var} [\bar{Y}]$ become easier for the approximation.

On the figure 2.8 we can compare the variance for different level of correlation.

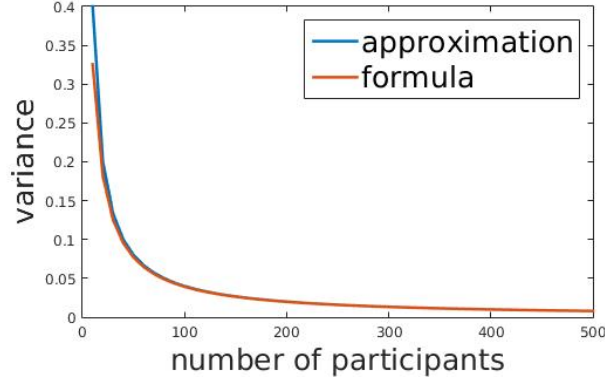


Figure 2.7: $\rho = 0.6$

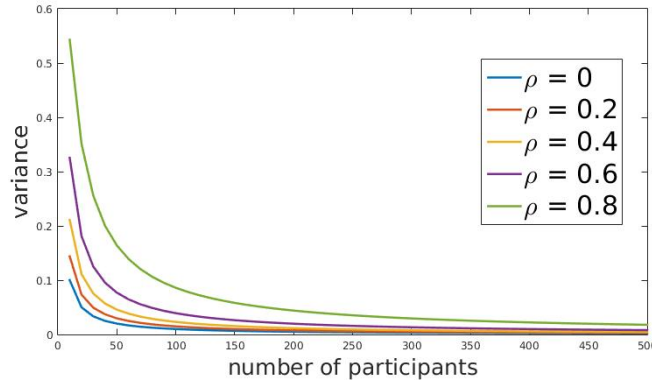


Figure 2.8: $\rho = 0.6$

2.4.3 Variance with skipping

Let's say that now we collected nk samples $Y_1, Y_2, Y_3, \dots, Y_{nk}$ that are correlated in the same way as in previous subsection. We will take each k sample and look at the variance of the next random variable:

$$\bar{Y}^k = \frac{Y_k + Y_{2k} + Y_{3k} + \dots + Y_{nk}}{n}$$

Let's note that the correlation between the variables Y_{ik} and $Y_{(i+h)k}$ is:

$$\text{corr}(Y_{ik}, Y_{(i+h)k}) = \rho^{kh}$$

If we introduce new random variables Z_1, Z_2, \dots, Z_n such that $Z_1 = Y_k, Z_2 = Y_{2k}, Z_3 = Y_{3k}, \dots, Z_n = Y_{nk}$ and $r = \rho^k, \bar{Z} = \bar{Y}^k$. Then:

$$\text{corr}(Z_i, Z_{i+h}) = \text{corr}(Y_{ik}, Y_{(i+h)k}) = \rho^{kh} = r^h$$

To sum up we have random variables Z_1, Z_2, \dots, Z_n where the correlation between any two of them depends of the distance $\text{corr}(Z_i, Z_{i+h}) = \rho^{kh} = r^h$. For this problem we already know the variance of \bar{Z} :

$$\text{var} [\bar{Z}] = \frac{\sigma^2}{n} \frac{1 - r^2 - 2r/n + 2r^{n+1}/n}{(1 - r)^2}$$

Or approximation:

$$\text{var} [\bar{Z}] \simeq \frac{\sigma^2}{n} \frac{1 + r}{1 - r}$$

Let's return to the previous notation and then we get:

$$\text{var} [\bar{Y}^k] = \frac{\sigma^2}{n} \frac{1 - \rho^{2k} - 2\rho^k/n + 2\rho^{k(n+1)}/n}{(1 - \rho^k)^2}$$

or:

$$\text{var} [\bar{Y}^k] \simeq \frac{\sigma^2}{n} \frac{1 + \rho^k}{1 - \rho^k}$$

2.4.4 In RDS context

Returning to the RDS context we can finally write the variance for different scenarios with the same budget. On the same budget B we have and taking each k node as a participant we will collect samples $Y_k, Y_{2k}, \dots, Y_{\frac{kB}{kC_1+C_2}}$, where number of participants is $m = \frac{B}{kC_1+C_2}$.

And if the samples are correlated as in the previous section then variance of the estimator, depending on k , will be:

$$\sigma_{\hat{\mu}_{SA}}^2(k) = \frac{\sigma^2}{\frac{B}{kC_1+C_2}} \frac{1 - \rho^{2k} - 2\rho^k/\frac{B}{kC_1+C_2} + 2\rho^{k(\frac{B}{kC_1+C_2}+1)}/\frac{B}{kC_1+C_2}}{(1 - \rho^k)^2}$$

Or approximated:

$$\sigma_{\hat{\mu}_{SA}}^2(k) = \frac{\sigma^2}{\frac{B}{kC_1+C_2}} \frac{1 + \rho^k}{1 - \rho^k} \quad (2.10)$$

We can regard this expression as at the function of k . Now in order to decide, how much samples to skip between the participant, we need to find the minimum of it.

On the figure 2.9 how the variance 2.10 changes for different level of dependency among values.

First let's observe what happens to the factors of the variance when we increase the k . The number of participants decreases, so the left factor in the expression 2.10 $\frac{\sigma^2}{\frac{B}{kC_1+C_2}}$ increases. In the same time the participants are less correlated, the right factor $\frac{1+\rho^k}{1-\rho^k}$ decreases. And the behavior of the variance depends on which factor is "stronger".

Let's for now concentrate on the case when $\rho = 0.8$. Starting when the $k = 1$ we observe that variance decreases when k increases. It means that correlation is very high and by skipping some values we reduce it significantly for the variance. When $k = 7$ we observe that function reaches its minimum. This result says that in these settings by taking only each seventh individual as a participant we will obtain the minimum error. As k further increase, the variance also start to increase slowly. It signals as that by skipping more than 7 persons, we will not decrease significantly the correlation, Trying to skip more than seven individuals we will just waste money on informants without any purpose.

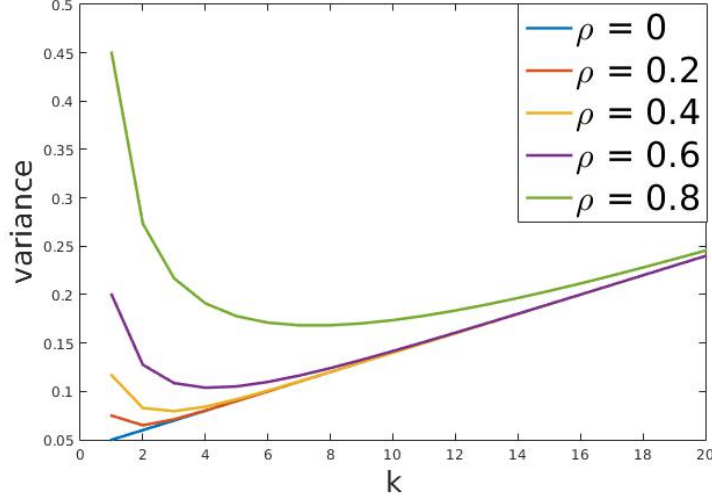


Figure 2.9: Variance with the formula 2.10 when $B = 100, C_1 = 1, C_2 = 4$

We have a trade-off here. If we skip too many values we can drastically decrease the number of participants. And the opposite, if we don't skip at all the correlation between values can be too high.

The value k when the variance is minimal depends on the level of correlation. Thus we observe that when we take lower values of ρ the desired value of k decreases. This coincides with our intuition: the lower is dependency, the less values we need to skip. Finally we see, that in case of no correlation ($\rho = 0$) it is useless to skip.

Moreover by studying the variance function, we can notice some interesting properties. Particularly, the observations that we stated in the section 1.5 were based on the observation and human logic. Now, we can show it formally.

Observation 1 *Just thinning of sample doesn't help*

Just thinning means that when all the samples are collected y_1, y_2, \dots, y_n we can try to take only k part of them. However to collect each sample we spend equal amount of money, let's denote it as $(C_1 + C_2)$. It means that we should look on function:

$$f(k) = \frac{\sigma^2}{B/(k(C_1 + C_2))} \frac{1 + \rho^k}{1 - \rho^k}$$

This function is strictly increasing when $k \geq 1$. It means that thinning the sample can only harm.
put the proof in the appendix

Observation 2 *Skipping reduces variance*

Here we have the fixed sample size n . However we can change the distance between samples as much as we want for free, $C_1 = 0$. Then we should look at the function:

$$f(k) = \frac{\sigma^2}{n} \frac{1 + \rho^k}{1 - \rho^k}$$

which is decreasing function.

This result is also very logical. When informants are not payed we should use them as much as we can.

2.4.5 Trying to use the result

Having the shown result we tried to use it. To use the function 2.9 to find the desired k which minimizes the error we need to define parameter ρ . What we tried is to take as ρ the average correlation between the immediate neighbors when the graph and the values on its nodes are given.

However, the suggested k did not coincide with the k found with the experimentations. The explanation could be following.

We that the values Y_1, Y_2, \dots, Y_n have correlation $\text{corr}(Y_i, Y_{i+h}) = \rho^h$. Correlation between values Y_1, Y_{101} is of the power 100. However when we walk on the graph the values Y_1 and Y_{101} can appear to be of the immediate neighbors. And when in reality the values of the corresponding nodes can be pretty much correlated, according to the geometrical model their correlation is so small that can be neglected.

The correlation between values is underestimated, except for the correlation of the two consecutive samples Y_i, Y_{i+1} . Therefore the correlation factor and its influence are underestimated. And when in reality skipping one more value would significantly reduce correlation, in the model it can be different.

put picture where is suggested k with this method does not coinc with exper results !!! the one with facebook

Therefore, the model is too simplified to use. In the next section we will generalize it on any graph with any values. However, knowing that the result from the model is underestimated, it can give us the *minimal* number of node to skip. And the result may not be the best possible, but comparing to the standard scenario has lower error.

2.4.6 General case

Due to the fact that the correlation between samples was oversimplified the results of geometric model were not exactly the same that in reality. This results could give us the lower bound on the value k , but we wanted to find it precisely.

In the previous example we were able to write the variance of the estimator because the correlation between all the random samples was known.

So our first goal is to generalize formula for the variance of the estimator when the values Y_1, Y_2, \dots, Y_n are collected with the random walk on the graph.

Let $f = f_1, f_2, \dots, f_n$ be the values of the attribute on the nodes $1, 2, \dots, n$. First, let P be the transition matrix of the random walk. We consider that probability for the individual to choose any of his friend is the same. Therefore, if the random walk visits the node i then one of the neighbors of this node will be chosen with probability d_i . Then we can write the transition matrix:

$$p_{ij} = \begin{cases} \frac{1}{d_i} & \text{if } i \text{ and } j \text{ are neighbors} \\ 0 & \text{if } i \text{ and } j \text{ are not neighbors} \\ 0 & \text{if } i = j \end{cases}$$

The stationary distribution of the random walk is:

$$\pi = \left(\frac{d_1}{\sum_{i=1}^n d_i}, \frac{d_2}{\sum_{i=1}^n d_i}, \dots, \frac{d_n}{\sum_{i=1}^n d_i} \right)$$

Let Π be the matrix that consist of n rows, where each row is the vector π . Let $f = (f_1, f_2, \dots, f_n)$ be the values of the attribute on the nodes $1, 2, \dots, n$.

We consider also that chain starts from initial distribution π . Then covariance between the random values Y_i and Y_j depends only on $j - i$:

$$\text{cov}(Y_i, Y_j) = \text{cov}(Y_0, Y_{j-i})$$

And then the formula for the variance:

$$\text{var}(Y_0) = \langle f, f \rangle_\pi - \langle f, \Pi f \rangle_\pi$$

and for the covariance:

$$\text{cov}(Y_0, Y_h) = \langle f, (P^h - \Pi)f \rangle_\pi$$

where $\langle a, b \rangle_c$ is weighted scalar product and if $a = (a_1, \dots, a_n), b = (b_1, \dots, b_n), c = (c_1, \dots, c_n)$ then:

$$\langle a, b \rangle_c = \sum_{i=1}^n a_i b_i c_i$$

To see, how these formulas were derived, consult 6 chapter of the book [7].

Using these formulas we can write the formula for the variance of the estimator as:

$$\begin{aligned} \text{var} [\bar{Y}] &= \frac{1}{n^2} \left(n \text{var}(X_i) + 2 \sum_{i=1}^n \sum_{j|i < j}^n \text{cov}(Y_i, Y_j) \right) = \\ &= \frac{1}{n^2} \left(n(\langle f, f \rangle_\pi - \langle f, \Pi f \rangle_\pi) + 2 \sum_{i=1}^n \sum_{j|i < j}^n \langle f, (P^{j-i} - \Pi)f \rangle \right) \end{aligned} \quad (2.11)$$

In this way we know the variance of the SA estimator, the expression is quite cumbersome. Another problem is that practically speaking computing the large powers of the matrix P can take a lot of time. Therefore, we will try to simplify this expression.

Let's look at auxiliary matrix P^* . Let D be the matrix $n \times n$ where $d_{ii} = \pi_i$ and $d_{ij} = 0$ if i is different from j . Auxiliary matrix P^* is build in the following way:

$$P^* = D^{\frac{1}{2}} P D^{-\frac{1}{2}}$$

Then if $\lambda_i (i = 1..r)$ are eigenvalues, v_i are corresponding right eigenvectors and u_i are corresponding left eigenvectors of the matrix P^* :

$$P^h - \Pi = \sum_{i=2}^r \lambda_i^h v_i u_i^T \quad (2.12)$$

For more explanation refer to the chapter 6 of the book [7].

Using formulas 2.11 and 2.12 and reasoning similar to the case with the geometric model we will derive following formula for the variance of the estimator SA:

$$\text{var} [\bar{Y}] = \frac{1}{n} \sum_{i=2}^r \frac{1 - \lambda_i^2 - 2 \frac{\lambda_i}{n} + 2 \frac{\lambda_i^{n+1}}{n}}{(1 - \lambda_i)^2} \langle f, v_i \rangle_\pi^2 \quad (2.13)$$

It means that for general model the correlation factor $f(n)$ is:

$$f(n) = \frac{1}{var(Y_0)} \sum_{i=2}^r \frac{1 - \lambda_i^2 - 2\frac{\lambda_i}{n} + 2\frac{\lambda_i^{n+1}}{n}}{(1 - \lambda_i)^2} < f, v_i >_{\pi}^2$$

Again as in the geometric model the expression 2.13 for the variance can be simplified if we take as the correlation factor when n approaches infinity:

$$var[\bar{Y}] = \frac{1}{n} \sum_{i=2}^r \frac{1 + \lambda_i}{1 - \lambda_i} < f, v_i >_{\pi}^2$$

Finally when we know the expression for estimator we can go back to the problem of sampling hidden population. Then for the different scenarios with different number of samples that are skipped between the participants $k - 1$ we get the variance:

$$var[\bar{Y}^k] = \frac{1}{\frac{B}{kC_1+C_2}} \sum_{i=2}^r \frac{1 - \lambda_i^{2k} - 2\frac{\lambda_i^k}{\frac{B}{kC_1+C_2}} + 2\frac{\lambda_i^{k(\frac{B}{kC_1+C_2}+1)}}{\frac{B}{kC_1+C_2}}}{(1 - \lambda_i)^{2k}} < f, v_i >_{\pi}^2 \quad (2.14)$$

or simplified version:

$$var[\bar{Y}^k] = \frac{1}{\frac{B}{kC_1+C_2}} \sum_{i=2}^r \frac{1 + \lambda_i^k}{1 - \lambda_i^k} < g, v_i >_{\pi}^2$$

Interestingly, the expression for the variance in general case has the same structure as for the geometric model. Therefore, all the proved observations for the geometric model are true for the general model. Moreover the interpretation of the derived formula is the same. There are two factors, left and right, that "compete" with each other. If we try to decrease the left factor, we will increase the right and the opposite. In order to find the desired parameter k we need to find the minimum of the estimator function for variance.

Even if it is difficult to obtain the explicit formula for k , the fact that k is integer allows us just to find it with search.

2.4.7 Error prediction

However variance it is not the only source of the error. We should also consider bias: the difference between expected value of the estimator and the real value. We can then compute bias of the estimator $\hat{\mu}_{avg}$ as following:

$$bias(\hat{\mu}_{SA}) = E[\hat{\mu}_{SA}] - \mu_{SA} < f, \pi > - \mu_{SA}$$

Then the mean squared error of the estimator, $MSE(\hat{\mu}_{SA})$, can be written with variance and bias as:

$$MSE(\hat{\mu}_{SA}) = bias(\hat{\mu}_{SA})^2 + var(\hat{\mu}_{SA})$$

We should note, that for all the scenarios with different k the bias is the same as this bias is due to the way that sampling is performed, random walk, that is the way of sampling for all the scenarios.

2.4.8 Discussion

The formula 2.14 implies that the graph and the values on the nodes are known. In this case we can correctly predict the error of the estimator. However, the researchers usually do not possess this information.

The expression 2.14 may be simplified for the graphs for which the distribution of eigenvalues is known. Like this it will be applicable, but it will lose generality.

There is no way to predict this value correctly without knowing the graph and the values of the nodes. The variance depends both on the graph structure. When researchers start respondent-driven sampling they may have no clue about graph structure. Therefore they can have only guesses about the error of the estimator that they will get.

2.5 Data

To validate our theoretical results we performed numerous simulations. For the graph structure with values on the nodes we used different sources of the data.

First, we used random graphs with values generated with the algorithm. Then we used real network structures like part of Facebook [6] and assigned values with the algorithm.

The real network structures where the nodes have some data are scarce. There are just several sources like this. We used data from the Project 90[5], data from the project Add health[4].

2.5.1 Data from the Project 90

Project 90 [5] studied how the network structure influences on the HIV prevalence. Besides the data about social connection the study collected some data about drug user, such as race, gender, whether he/she is sex worker, pimp, sex work client, drug dealer, drug cook, thief, retired, housewife, disabled, unemployed, homeless.

For our experiments we took the largest connect component from the available data, which consists from 4430 nodes and 18407 edges.

On the figure 2.10 you can see the graph structure built with the Gephi tool [2], where the attribute gender is depicted with color for every node.

2.5.2 Data from the Add health project

The National Longitudinal Study of Adolescent to Adult Health (Add Health) is a huge study that began with surveying students from the 7-12 grades in the United States during 1994-1995 school year. In general 90,118 students representing 84 communities took part in this study. The study kept on going survey student as they growing up. The data includes information about social, economic, psychological and physical status of the students and other.

The social network of students' connections was built based on the reported friends. Each of the students was asked to provide the names of 0-5 male and 0-5 female friends. Then the network structure was built that now can help to analyze if some characteristics of the students indeed are influenced by their friends.

Though this data are very valuable they are not in the free access. Part of them are actually available but the ids of the student and of his/her friends are masked that makes impossible to recreate the network. However the part of the data can be accessed through the link [3] but only

with few attributes for students, such as: sex, race, grade in school and, if communities that have two schools, the school code (explain better about sister school).

There are several network for different communities. On the figure 2.11 represented the graph with 1996 nodes and 8522 edges, built again with the Gephi tool [2], where the attribute race is depicted with color for every node.

2.5.3 Simulated form Gibbs distribution

2.5.4 Experiments

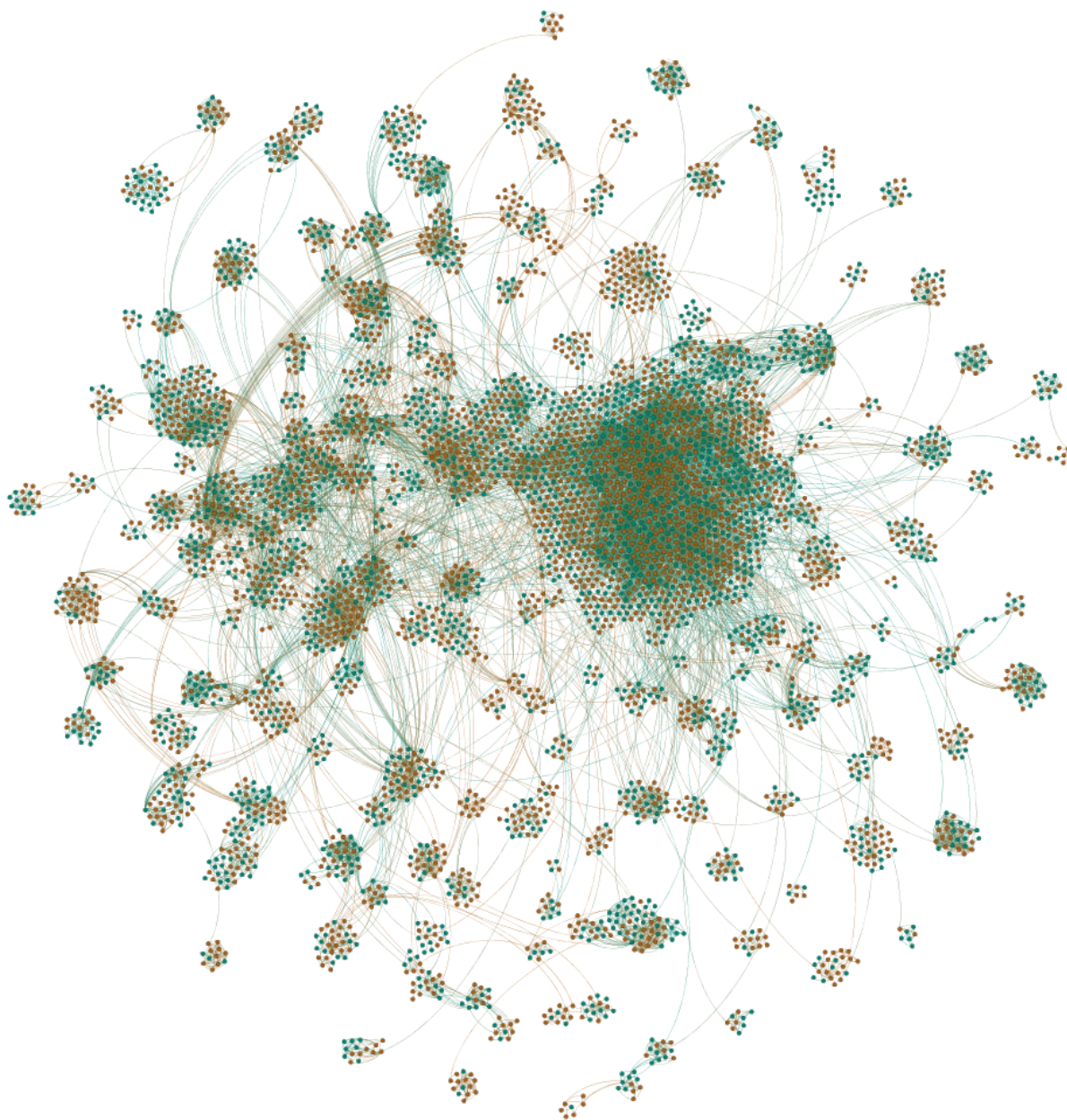


Figure 2.10: Graph from the Project 90 data. The colors of the node represent the gender: brown nodes correspond to the male participant, green nodes correspond to the female participants

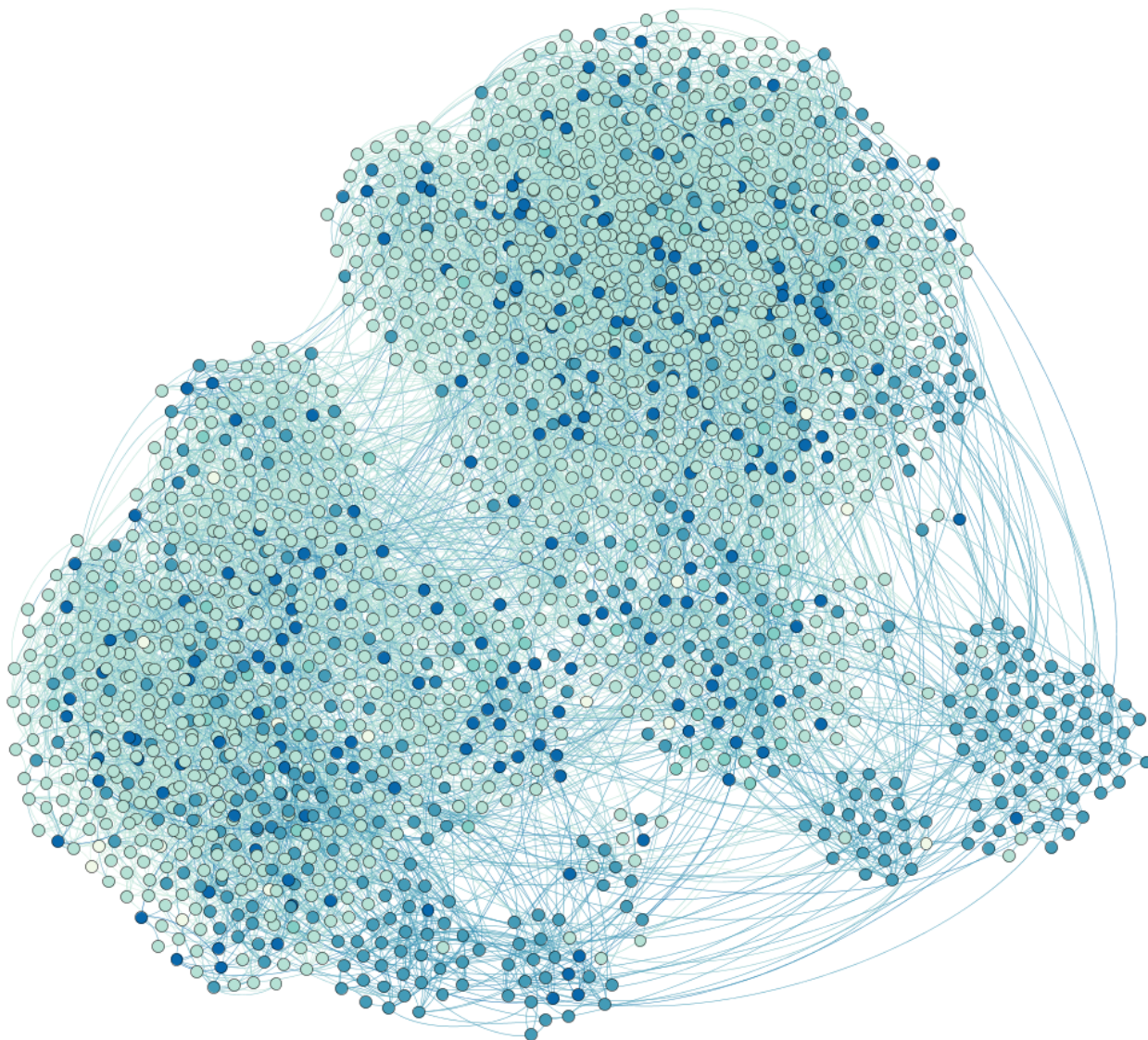


Figure 2.11: Graph from the Add Health data. The colors of the node represent the race

Chapter 3

Comparison to other methods

Estimator 1

Estimator 2

Estimator 3

As stated also [12] the advantage to use estimator1 appears only when the needed for estimation trait depends on the degree of the node. To support this statement they compare particularly the standard error of the sample mean and RDS estimate on the data sets from Project 90 and Add Health. The results are presented on the figure [put pictures].

3.1 Further study

Conclusion

In this work we regarded the sampling technique that is called respondent-driven sampling. This technique can be very useful in the cases when the members of the population are hard to find. The fact that members of the population form a network allows researcher to use their connections to reach another members.

We observed that this way of sampling is not uniform independent sampling. The way of sampling and the presence of homophily in the network influence on the error of the estimator. The researchers should keep attention that particularly the variance of the estimator is grater than in the uniform independent sampling.

The level of the homoplily in the network influences on the correlation between samples. One way of decreasing the correlation is thinning out the sample. We showed formally that just thinning the sample will increase the error of the estimator. Instead we proposed to enhanced respondent-driven sampling that allows to decrease the correlation between samples without reducing drastically sample size.

In the same way that random graphs can imitate the structure of the graph with needed correlation we needed the mechanism of assigning values to the nodes in such a way that it imitates the property of homophily in the network. Created algorithm allows not only to do this but also to control the level of correlation in the network. This result is not applicable to this particular model. But it is the general mechanism that allows to create the network with controllable level of dependency between neighbours. That is highly valuable for research purposes.

Using created mathematical model we targeted to find the error of the estimator for different scenarios of enhanced RDS. The scenario with the minimal error would be the best and recommended for applying. First we regarded the geometric model where correlation between samples is the power of difference between observing time of two samples.

Though this model turned out to be too simplify to produce the exact result it gave us valuable products. First, some . Second, it can give us the lower bound.

Then we were able to generalize the result for any kind of graph with any kind of correlation. Theoretical results were validated with experimentations. For the network with values with used created model and the real networks as well.

to correct: - cautions notes - change n on m in the enhanced RDS - put the transition matrix in the beginning of the second chapter - everywhere expected energy of random field not configuration!!! restructure report appeared - j turned out

geometric model allow to show some rasults taht are easy generalized to the general case
search everywhere: instead variation - variances

Bibliography

- [1] Computer science university of maryland. <http://www.cs.umd.edu/hcil/science20>. Accessed: 2015-08-07.
- [2] Gephi - the open graph viz platform. <http://gephi.github.io/>.
- [3] Linton c. freeman, research professor, department of sociology and institute for mathematical behavioral sciences school of social sciences, university of california, irvine. <http://moreno.ss.uci.edu/data.html>. Accessed: 2015-07-01.
- [4] The national longitudinal study of adolescent to adult health. <http://www.cpc.unc.edu/projects/addhealth>. Accessed: 2015-07-01.
- [5] The office of population research at princeton university. <https://opr.princeton.edu/archive/p90/>. Accessed: 2015-07-01.
- [6] Stanford large network dataset collection. <https://snap.stanford.edu/data/>. Accessed: 2015-07-01.
- [7] Pierre Brémaud. *Markov chains: Gibbs fields, Monte Carlo simulation, and queues*, volume 31. Springer Science & Business Media, 2013.
- [8] Nicholas A Christakis and James H Fowler. The spread of obesity in a large social network over 32 years. *New England journal of medicine*, 357(4):370–379, 2007.
- [9] Krista J Gile and Mark S Handcock. Respondent-driven sampling: An assessment of current methodology. *Sociological methodology*, 40(1):285–327, 2010.
- [10] Minas Gjoka, Maciej Kurant, Carter T Butts, and Athina Markopoulou. A walk in facebook: Uniform sampling of users in online social networks. *arXiv preprint arXiv:0906.0060*, 2009.
- [11] Sharad Goel and Matthew J Salganik. Respondent-driven sampling as markov chain monte carlo. *Statistics in medicine*, 28(17):2202–2229, 2009.
- [12] Sharad Goel and Matthew J Salganik. Assessing respondent-driven sampling. *Proceedings of the National Academy of Sciences*, 107(15):6743–6747, 2010.
- [13] Douglas D Heckathorn and Joan Jeffri. Jazz networks: Using respondent-driven sampling to study stratification in two jazz musician communities. In *Unpublished paper presented at American Sociological Association Annual Meeting*, 2003.

- [14] Bruno Kauffmann, François Baccelli, Augustin Chaintreau, Vivek Mhatre, Konstantina Papa-
giannaki, and Christophe Diot. Measurement-based self organization of interfering 802.11 wire-
less access networks. In *INFOCOM 2007. 26th IEEE International Conference on Computer
Communications. IEEE*, pages 1451–1459. IEEE, 2007.
- [15] Helgar Musyoki, Timothy A Kellogg, Scott Geibel, Nicholas Muraguri, Jerry Okal, Waimar Tun,
H Fisher Raymond, Sufia Dadabhai, Meredith Sheehy, and Andrea A Kim. Prevalence of hiv,
sexually transmitted infections, and risk behaviours among female sex workers in nairobi, kenya:
Results of a respondent driven sampling study. *AIDS and Behavior*, 19(1):46–58, 2015.
- [16] Michael Pollard, Harold D Green, David P Kennedy, Myong-Hyun Go, and Joan S Tucker.
Adolescent friendship networks and trajectories of binge drinking. 2013.
- [17] Jesus Ramirez-Valles, Douglas D Heckathorn, Raquel Vázquez, Rafael M Diaz, and Richard T
Campbell. From networks to populations: the development and application of respondent-driven
sampling among idus and latino gay men. *AIDS and Behavior*, 9(4):387–402, 2005.
- [18] Matthew J Salganik and Douglas D Heckathorn. Sampling and estimation in hidden populations
using respondent-driven sampling. *Sociological methodology*, 34(1):193–240, 2004.
- [19] Parag Singla and Matthew Richardson. Yes, there is a correlation:-from social networks to
personal behavior on the web. In *Proceedings of the 17th international conference on World
Wide Web*, pages 655–664. ACM, 2008.