# Book of Short Papers
# SIS 2020

SIS2020 *Pisa*

SIS
Società
Italiana di
Statistica

Editors: Alessio Pollice, Nicola Salvati and Francesco Schirripa Spagnolo

# Contents

## Specialized sessions

# Solicited Sessions

# Contributed papers and Posters

# Environment, Physics and Engineering ......................................................809

## Finance, business and official statistics ......................................886

# Machine Learning and Data Science.......................................................1023

## Models and methods - Categorical, Ordinal, Rank Data .........................1135

## Models and methods – Regression.........................................................1206

# Models and methods – Sampling ...............................................1271

# Models and methods - Theoretical Issues in Statistical Inference ..........1314

# Models and methods - Time Series and Longitudinal Data.....................1350

## Population and society .................................................................................1411

# Evaluating a Hybrid One-Staged Snowball Sampling through Bootstrap Method on a Simulated Population

## Campionamento a valanga ibrido ad uno stadio: una valutazione con metodo bootstrap su una popolazione simulata

Venera Tomaselli[1], Giulio Giacomo Cantone[2]

**Abstract:** Snowball sampling is a design where a number of individuals is surveyed in a population and then is requested to share the survey tool among their 'social links'. The aim is to recruit other people into the sample. One-stage recruitment means that individuals who are recruited by people already participant are not requested to recruit any others. Although surveys adopting snowball are financially less expensive than alternatives, the reliability of asymptotical estimates through this method is often questioned. A hybrid snowball sampling is designed such that a quota is randomly sampled, and another through snowball. Through a simulation of a demographic setup we found that bootstrap statistics of a one-staged hybrid design asymptotically show no significant difference to random design. This does not extend into situations where the random quota is smaller. We conclude that more complex setups will be needed before to generalise these results.

**Abstract:** *Il campionamento a valanga è un disegno campionario dove un sondaggio è sottoposto ad certo numero di individui, ed è poi chiesto a questi di condividere il sondaggio tra i loro 'legami sociali' al fine di reclutarne altri nel campione. Il reclutamento ad un livello prevede che chi è reclutato da chi è già partecipante non ne recluti altri. Sebbene i sondaggi a valanga siano meno costosi delle alternative, l'affidabilità delle stime di questo metodo è spesso posta in dubbio. Un campionamento a valanga ibrido è un disegno campionario dove una frazione del campione è estratta a sorte ed un'altra reclutata a valanga. Per mezzo di una simulazione di un setup demografico abbiamo trovato che le statistiche bootstrap di un disegno ibrido ad uno stadio non mostrano differenze significative rispetto all'alternativa casuale. Ciò non si estende ad una situazione con una minor frazione di individui estratti a caso. Riteniamo che setup più complessi siano necessari prima di generalizzare questi risultati.*

**Key words:** hybrid one-staged snowball sampling, bootstrapping, simulated population.

## 1 Introduction to snowball sampling

The terminology 'snowball sampling' denotes the methodology of social research where "a small [random] sample of persons [is interviewed], asking who their best friends are, interviewing these friends, then asking them their friends, interviewing these, and so on" (Coleman, 1958, 29). Leo Goodman (1961), who was a colleague of Coleman at University of Chicago, adopted "snowball sampling" to refer to a

---

[1] Venera Tomaselli, Department of Political and Social Sciences, University of Catania IT, e-mail: venera.tomaselli@unict.it (*corresponding author*).

[2] Giulio Giacomo Cantone, Department of Physics and Astronomy, University of Catania IT, e-mail: prgcan@gmail.com.

mathematical model to "make statistical inferences about various aspects of the relationships present in the population" by "data obtained using an s stage k name snowball sampling procedure" (*idem*, 148).

The aim of Goodman was to formalize conditions such to assume asymptotically unbiased estimates from Coleman's procedure. Goodman's model assumes that $k$, the amount of interviewed people per $s$ stages of the chain, is a constant. "Stage 0" or $s = 0$ is the initial sample of $k$ cases in the population and is random, instead. Snijders (1992) offered a different formal model with no fixed $k$. Both authors conclude that while a snowball procedure is not a fully randomized method, if 'homophily' among participants is sufficiently low, the sample may be assumed to be representative of a population.

Homophily is a terminology introduced since the '50 in sociological works and later formalized by network scientists as a parametrization of the correlation between the value of the variable in a node and the amount of its edges (Newman, 2010). Nevertheless, the interpretation of homophily is often ambiguous. Crawford, Aronow, Zeng et al. (2017) claim that homophily is often confused with "preferential recruitment" which is when the probability to be drawn in the sample from stage 1 cannot be assumed uniformly distributed.

When Goodman (2011) returned on the topic of snowball sampling, he highlighted an issue: while the model he developed assumed a random primary sample ('stage 0'), many studies performed after lacked this assumption[3]. However, Craig, Hays, Pickard et al. (2013) found that on 7 surveyed panel vendors, 6 had a mean of 20% of their proposed participants involved in at least one other of the 6. This leads to intuitively think that samples can support a certain 'quota' of biased cases within, without a great loss of precision in estimates within a multivariate design.

The research of Etter and Perneger (2000) is very noteworthy in this sense. Authors surveyed a random sample of 1000 residents in Geneva aged 18-70 (primary participants) in 1997. They asked every contacted subject, even those not willing to participate, to mail the questionnaire to all the smoker and ex-smoker residents in Geneva they knew (secondary participants). At the end of the data gathering process, 3,300 residents were mailed with the questionnaire and 1,167 individuals (35%) returned the questionnaire filled. Of these, 578 were primary participants and 566 were secondary participants. According to authors the mean age difference between the samples was only of 1.7 years ($p$-value=0.003). The only other significant difference was in sex ratio (7% difference, $p$-value=0.009) while behavioural traits were reported to be not significantly different in the two groups. As unintended consequence, the authors obtained two samples very similar in size within the one-staged snowball. We call this sampling design 'hybrid one-staged snowball sampling'.

In a study comparing a random statistical approach to a snowball-based approach to estimate morbidity and mortality of the rare disease visceral Leishmaniosis in two districts of the state of Bihar, India in 2011, Siddiqui, Rabidas, Sinha et al. (2016) concluded that snowball approach was found not sensitive enough to be adopted to estimate morbidity of the rare disease. At the same time, authors noticed that comparing costs, snowball approach required 1/6 of man-days and half the financial costs of the random alternative. Since full snowball design could not be suited for many applications, we want to evaluate if adoption of hybrid sampling can overcome losses of accuracy while reducing costs of research.

## 2 Simulation of a hybrid sampling

Testing procedures on empirical data in population studies raises the following issues:
- usually empirical research is designed for finding information over a phenomenon, not to provide data for evaluative research about methods
- we lack 'true values' on parameters of variables. We can test the hypothesis if two samples are drawn from the same population within confidence level (CL) but this says nothing about which design is better performing on estimation.

We propose instead a combination of two computational methods:

---

[3] According to the author, this was a result of historical label of 'snowball' mostly for the Coleman's procedure of data collection. Hence, the technique (also referred as 'chain-referral sampling') was associated to qualitative studies on 'rare population' or on 'hidden-populations' (Atkinson and Flint, 2001). Erikson (1979) and Snijders (1992) expressed concern about the actual possibility to randomly sample the stage 0 for hidden populations.

- a simulation model to procedurally generate a virtual dataset of population data from known parameters. Simulated populations can be modeled to fit pre-existing data with a small error (Alfons, Kraft, Templ et al., 2011; Burgard, Kolb, Murkle et al., 2017) or according with a theoretical model i.e. to construct hypotheses
- computational *bootstrap* sampling, falling under the more generic terminology of Monte Carlo methodologies (Gil, Montenegro, González-Rodríguez et al., 2006; Gobet, 2016).

We want to test performance of a sampling design against a standard. Hence, the simulated population does not need to fit empirical data. Still, this virtual population needs to be generated according to non-unrealistic assumptions over expected general outcomes. In particular, assumptions will take form of structural equations made through random variables. Technical sophistications in assumptions of the model may be added once the validity of the simplest models is established.

Bootstrap means that the model is then run many times, and each time estimates are recorded. Hypotheses are tested on the statistics (e.g. the so-called bootstrap-mean) of the distribution of values of the recorded estimates.

The general issues we found for the proposed simulation are two:

*Issue 1*: How can we simulate non-unrealistic variables in a virtual population?

In order to reflect the complexity of population studies without losing simplicity of linear structural equations, N simulated individual in the population (or 'agent' of the model) will be designed such to show few random variables:
- age: $X_0$; Gompertz and Weibull are established density functions to model age (Ricklefs and Scheuerlein, 2002). As the Weibull's function has only two parameters, *scale* ($\lambda$) and *slope* (k), we adopt it because it is a simpler assumption for our purposes[4].We propose:

$$X_0 \sim Weibull\left(\lambda_{X_0} = 50, k_{X_0} = 3\right).$$ [1]

For ethical issues, data on minors are usually not collected; therefore, any agent with $X_0 < 18$ (~ 4,56% of N) will be removed from the virtual population.
- sex: $X_1$; a dummy variable with a probability equal to 0.5 per side:

$$p(X_1 = 1) = 0.5$$ [2]

- behaviour: $X_2$; the research on behavioural metrics is a vast field with many insightful contributions on proper tools to measure behaviour's presence and extension (Kline, 2015). But we aim at keeping unambiguity and simplicity in assumptions before prompting actual 'field knowledge' into it. Therefore, we will just notice the presence/absence of the behaviour into a dummy variable with a probability density. For the proposal of the simulation, we arbitrarily set again 0.5 per side:

$$p(X_2 = 1) = 0.5$$ [3]

- chronic condition: $Y$; the outcome of a structural equation where $X_0$ of the agent is the baseline, $X_1$ and $X_2$ are endogenous factors, $\beta$ are weights, and $\varepsilon$ is a continuous random error variable which sums all exogenous factors of $Y$:

$$Y = \begin{cases} 0, X_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon < t_\zeta \\ 1, X_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon > t_\zeta \end{cases}$$ [4]

$t$ is an integer threshold value such that the probability of $Y=0$ in the whole population asymptotically tends to $\zeta$. To intuitively understand the equation: if we keep out the agents when $X_0 < 18$, set $\beta$ to be equals to 0 and $\varepsilon$ is symmetrical with an expected value equal to 0, then, for $\zeta = 0.5$, the value 45 approximates $t_{\zeta =0.5}$. In simple terms, 45 will be the median age of the simulated population, given these parameters.

We can use this information to model our assumptions on $\beta$. If we assume:
- normal distribution of $\varepsilon$ with $\mu=0$
- age, sex, and behaviour having all equal impact on chronic condition
- $\beta_1 = \beta_2 = 90$

then $t_{\zeta=0.5} = 135$[5].

The threshold $t$ can be set equal to 0 by the property of linear transformation of normal distributions (Bryc, 1995). Once $-(t_\zeta)$ is added to $\mu_\varepsilon$, we notice that the lower the $\mu_\varepsilon$, the higher the $\zeta$. Therefore, we

---

[4] The Gompertz's assumptions about the parameters may be more suited for fitting synthetic data, instead.

[5] $E(x_0) + E(\beta_1 x_1) + E(\beta_1 x_2) = t_{\zeta=0.5}$.

computationally find a setup pair of values for $\mu_\varepsilon$ and $\sigma_\varepsilon$ such that the expected value of $\zeta$ is approximates a desired value. We want a model where the chronic condition afflicts 1 agent every 20 in the population: for a population of 100,000 agents, decreased to ~ 95445 agents after removal of agents aged < 18, under $\zeta \sim 0.95$, we expect ~ 4775 agents afflicted by the chronic condition. In the next setup, the $E(Y) = 0.0505$

$$Y = \begin{cases} 0, X_0 + 45(X_1) + 90(X_2) + \varepsilon < 0 \\ 1, X_0 + 45(X_1) + 90(X_2) + \varepsilon > 0 \end{cases} \quad [5]$$

$$\varepsilon \sim Gaussian(\mu_\varepsilon = -295, \sigma_\varepsilon = 90)$$

$\sigma_\varepsilon$ was arbitrarily set such that $\sigma \sim 6(\sigma_{X_0})$.

*Issue 2*: How can we simulate recruitment in the hybrid sampling?

There are two different processes to simulate a network among individuals of our population. The first is to actually employ a software that connects all the members of the virtual population in a network. The 'primary' sample is randomly drawn, then agents in the primary sample recruits 1 stage of secondary participants among their links. Through this method, both homophily clusters and preferential recruitment are simulated. Unfortunately, this method is unpractical for a high-sized N of population as it requires intense computational resources.

Another approach to simulate hybrid sampling is less computationally intensive:

1. the primary sample $I$ is drawn randomly and removed from the $N$ population. Then two new variables are assigned to any $i$ agent that is an element of $I$[6]:
   - Links:
   $$Z_1 \sim Integer\left(ChiSquared\left(df_{Z_1} = 150.5\right)\right) \quad [6]$$
   - Recruitment:
   $$Z_2 \sim Integer\left(ChiSquared\left(df_{Z_2} = 5.5\right)\right)^{[7]} \quad [7]$$
2. every $i$ agent in the first sample randomly draws $z_1$ other agents $j_i$ from the $N$ population; $j_i$ agents are not removed from $N - I$ but becomes elements of $J_i$
3. for each $i$ with a $z_2 > 0$ and all its $j_i$, a $d_{(i,j)}$ value is assigned. $d_{(i,j)}$ measures statistical distance between $i$ and $j$ expressed by the value of a structural equation

$$D_{(i,j)}(X_0, X_2): d_{(i,j)} = \left|x_{0_i} - x_{0_j}\right| + \varepsilon_d \quad [8]$$

$$\varepsilon_d \sim \left|\left(Gaussian(\mu_\varepsilon = 0, \sigma_\varepsilon \sim 0)\right) - \left(Gaussian(\mu_\varepsilon = 0, \sigma_\varepsilon \sim 0)\right)\right| \simeq 0^{[8]}.$$

For each $i$, $z_2$ elements of $J_i$ are drawn with a probability equal to:

$$1 - \frac{d_{(i,j)}}{\sum d_{j_i}} \quad [9]$$

which is the function of preferential recruitment. In order to reduce endogenous bias introduced by how the structural equation of distance $D_{(i,j)}$ is modelled we perform a re-sample: a sample equal in size of $I$ is randomly re-sampled into $I'$ from the $\cup(J)$ union of all the $j$ agents drawn from each $i$. The union of $I$ and $I'$ is the actual 'hybrid sample'.

In this model only age has an impact on statistical distance. Impact of sex and other traits was already a topic in Coleman (1958), but since $E(X_1)$ and $E(X_2)$ are equal per side, any non-complex interaction among endogenous variables in $D_{(i,j)}$ would result in a more 'randomised' estimates in the final sample.

## 3 Results and conclusions

The simulation and all the procedures were performed with software $R^{[9]}$.

---

[6] This design reflects "layers of friendships" mentioned in Mac Carron, Kaski, and Dunbar (2016). The parameters are then chosen only as standards for the recruitment process.

[7] *Integer* is a function that subtracts to a value its mantissa. This operation was made to force integer values on the random variables. 0.5 was added to the parameters in order to compensate the loss of the mantissa. A different approach is to model the random variables from Poisson's function, or from an exponential one.

[8] While in this model the distance error $\varepsilon_d$ representing exogenous factors of $D_{(i,j)}$ is assumed to be equal to 0 by forcing $\sigma_\varepsilon \sim 0$, we kept it in the equation the normal model that in our opinion best fit a random error in distance for future developments.

[9] Packages: base, survival, dplyr, rlist.

The population of N ~ 95445 was simulated only once. The random sample was drawn with a size equal to 1056 agents, which is the minimum size of representative random sample for confidence interval of 3% and a CL equal to 95%. Hybrid sample was drawn from the same population but with a size of $I$ equal to 528, so the size of the final sample $I \cup I'$ is again 1056.

Both the design were run and recorded 300 times. We noticed that the average amount of $J \sim 2380$. Hence, we decided to perform a second hybrid sampling with $I = 264$ and $I' = 792$. Statistics are reported in table 1 and table 2.

**Table 1**. Population statistics.

| | $N$ | $X_0$ | | $X_1$ | | $X_2$ | | $Y$ | |
|---|---|---|---|---|---|---|---|---|---|
| | | $\mu$ | $\sigma^2$ | $\mu$ | $\sigma^2$ | $\mu$ | $\sigma^2$ | $\mu$ | $\sigma^2$ |
| Simulated Population | 95443 | 46.12 | 284.86 | 0.5 | 0.25 | 0.5 | 0.25 | 0.05 | 0.047 |

**Table 2**. Bootstrap statistics (n=300) of samples in random and hybrid sampling design.

| | $n$ | $X_0$ | | $X_1$ | | $X_2$ | | $Y$ | |
|---|---|---|---|---|---|---|---|---|---|
| | | Mean | Var | Mean | Var | Mean | Var | Mean | Var |
| Random Sampling | 300 | 46.15 | 0.22 | 0.49 | ~0 | 0.49 | ~0 | 0.049 | ~0 |
| Hybrid Sampling (stage 0 = 528) | 300 | 45.95 | 0.21 | 0.5 | ~0 | 0.49 | ~0 | 0.05 | ~0 |
| Hybrid Sampling (stage 0 = 264) | 300 | 47.5 | 0.25 | 0.5 | ~0 | 0.49 | ~0 | 0.051 | ~0 |

We checked if the difference between the mean of population (46.12) and the bootstrap means of all the variables in the random and hybrid samplings (see Table 2) are significant with $p$-value $< 0.05$. We found:
- no significant difference in all variables in random design, as expected
- no significant differences in variables of hybrid design (stage 0 = 528)
- significant difference in $X_0$ and Y in hybrid design (stage 0 = 264).

For this reason, we exclude feasibility of hybrid sampling (stage 0 = 264). On the basis of the results from our experiment, we can assert that exists at least a setup of very simple parameters under which hybrid model can be employed instead of random sampling design.

We think the variances in bootstrap statistics are abysmal because complexity in population data is low and recruitment processes show a low variance (i.e. Chi Square functions are leptokurtic).

# 4 Future developments

We propose the following topics as future developments of the present study:
1. Complexity in the population's equations: in empirical demographic data, age is not independent from sex. Behaviours are complex: some may be dependent from age and sex while others may not. For future applicative applications, the target behaviour may be modelled as a structural equation. A dummy model may result simplistic, so the output may be rescaled into a multipoint scale, instead. Chronic conditions may be results of non-linear interactions among factors. The final model could be driven from a scientific theory. For internal validity of the model, the relevant parameter is the skewness of the variables in the population: if the variables are not symmetrical, even in a population with zero homophily, where $J_i$ are randomly sampled, and if preferential recruitment is modelled such that agents in the $I$ subsample 'prefer' to draw other agents that are not statistically distant from them, then we expect an increase in estimates' variance.
2. Complexity in the preferential recruitment: we suppose this point is both the most controversial and the most impactful on biases in asymptotical estimates. We already stated that in order to simulate homophily, we would need to simulate a network. We think a practical compromise could

be to import a network dataset and then test sampling through it, although this incurs in issue mentioned in section 2. We admit we don't know how people recruit other respondents into survey tool, neither we feel like we can generalise too much on this issue. Our general intuitions for future developments are: (i) $Z_2$ should be platykurtic and with a fat right tail, i.e. an exponential distribution. The fatter the right tail, the more the model is stressed; (ii) the weight in likelihood of recruitment should be positively correlated to a statistical distance.

3. Multi-stage and stage-free models: in order to stress more a reduction in size of 'stage 0' *I*, more stages of recruitment can be added, so that, i.e. a two-staged hybrid sampling would sample the union of *I, I', I''*, etc., where the latter stage is recursively not randomly drawn among 'friends' of the precedent stage. If a development of the model simulates a network population with complex variables, a more realistic approach may be free of stages: in other terms, every agent can recruit another, until the union of all sampled agents is equal to target *n*.

4. To evaluate with a regressive model the fitness of bootstrap statistics: in the present study, the validity of the sampling design was tested through significance of difference between bootstrap statistics and population parameters. A different evaluative approach is to estimate fitness of bootstrap statistics in a regressive model for each Y, already known as outcome of structural equation.

# References

1. Alfons, A., Kraft, S., Templ, M., & Filzmoser, P. (2011). Simulation of close-to-reality population data for household surveys with application to EU-SILC. Stat Method Appl, 20(3), 383–407.
2. Atkinson, R. & Flint, J. (2001). Accessing Hidden and Hard-to-Reach Populations: Snowball Research Strategies. Social Research Update, 33.
3. Bryc, W. (1995). The Normal Distribution: Characterizations with Applications. Springer-Verlag.
4. Burgard, J. P., Kolb, J.-P., Merkle, H., & Münnich, R. (2017). Synthetic data for open and reproducible methodological research in social sciences and official statistics. AStA Wirtschafts- Und Sozialstatisches Archiv, 11(3-4), 233–244.
5. Coleman, J. (1958). Relational Analysis: The Study of Social Organizations with Survey Methods. Hum Organ, 17(4), 28–36.
6. Craig, B. M., Hays, R. D., Pickard, S. A., Cella, D., Revicki, D. A., & Reeve, B. B. (2013). Comparison of US Panel Vendors for Online Surveys, J Med Inter Res, 15(11).
7. Crawford, F. W., Aronow, P. M., Zeng, L., & Li, J. (2017). Identification of Homophily and Preferential Recruitment in Respondent-Driven Sampling. Am J Epidemiol, 187(1), 153–160.
8. Erickson, B. H. (1979). Some Problems of Inference from Chain Data. Sociol Methodol, 10, 276.
9. Etter, J.F., & Perneger, T. V. (2000). Snowball sampling by mail: application to a survey of smokers in the general population. Int J Epidemiol, 29(1), 43-48.
10. Gil, M. Á., Montenegro, M., González-Rodríguez, G., Colubi, A., & Rosa Casals, M. (2006). Bootstrap approach to the multi-sample test of means with imprecise data. Comput Stat Data An, 51(1), 148–162.
11. Gobet, E. (2016). Monte-Carlo Methods and Stochastic Processes. New York: Chapman and Hall/CRC.
12. Goodman, L.A. (1961). Snowball sampling. Annals of Mathematical Statistics. 32(1), 148–170.
13. Goodman, L.A. (2011). Comment: On Respondent-Driven Sampling and Snowball Sampling in Hard-to-Reach Populations and Snowball Sampling Not in Hard-to-Reach Populations. Sociol Methodol, 41(1), 347–353.
14. Kline, P. (2015). A Handbook of Test Construction (Psychology Revivals). London: Routledge.
15. Mac Carron, P., Kaski, K., & Dunbar, R. (2016). Calling Dunbar's numbers. Soc Netw, 47, 151–155.
16. Newman M (2010). Networks: An Introduction. New York, NY: Oxford University Press.
17. Ricklefs, R. &Scheuerlein, A. (2002). Biological implications of the Weibull and Gompertz models of aging. J Gerontol A-Biol, 57(2), 69-76.
18. Snijders, T. A. B. (1992). Estimation on the Basis of Snowball Samples: How to Weight? BSM, 36(1), 59–70.
19. Siddiqui, N. A., Rabidas, V. N., Sinha S. K., Verma R. B., Pandey, K. P., Singh, V.P., Ranjan, A., Topno, R. K., Lal, C. S., Kumar, V., Sahoo, G.C., Sridhar, S., Pandey, A., Das, P. (2016) Snowball Vs. House-to-House Technique for Measuring Annual Incidence of Kala-azar in the Higher Endemic Blocks of Bihar, India: A Comparison. PLoS Neglected Tropical Diseases, 10(9).