

Undergraduate Students Can do Meaningful Research

Sat Gupta¹, Igor V. Erovenko¹, Hyunju Oh³, Jan Rychtář¹ and Dewey Taylor²

¹*Department of Mathematics and Statistics, UNC Greensboro, NC, USA*

²*Department of Mathematics, Virginia Commonwealth University, Richmond, VA, USA*

³*Department of Mathematics and Computer Science, Bennett College, NC, USA*

Final Version Received on 23.08.2017

Abstract

This paper expands on some results presented by the first author at the annual conference of the Society of Statistics, Computer and Applications held in Jammu, India on March 6-8, 2017. We highlight the significance of undergraduate research and share examples of research projects undergraduate students are capable of handling. Our aim is to dispel the myth that undergraduate students cannot do meaningful research. With proper mentoring, undergraduate students can be successfully motivated to work on open research problems, make substantial progress towards a solution, present their work at professional conferences, and publish their work in peer-reviewed journals. We also discuss the liberal support undergraduate research receives in the United States from such funding agencies as the National Science Foundation, Mathematical Association of America, and the National Security Agency. It is our hope that other senior researchers will be inspired by our experience and will consider mentoring undergraduate student researchers.

Key Words: RRT Models, Greenberg Unrelated Question RRT Model, Game Theory, Sensitivity Analysis

1. Introduction

Numerous studies indicate that undergraduate research experiences lead to improved confidence and creativity, which in turn contribute to an increased likelihood to pursue graduate education (Russell, 2007; Lopatto, 2007; and Seymour, 2004). Lopatto (2009) demonstrates the significance of undergraduate research in attracting and retaining talented students in science. Moreover, undergraduate research enables underrepresented groups to succeed in scientific careers.

Undergraduate research receives significant attention in the United States. Even talented high school students are getting involved (see Gargeya and Pratap, 2015, and Goel and Tashakkori, 2015). Yet, many people still believe that undergraduate students are incapable of doing meaningful research. This paper intends to dispel that myth. In our experience, such belief is more prevalent outside the United States.

Major funding agencies, such as the National Science Foundation, Mathematical Association of America, and the National Security Agency, offer liberal funding under major programs like REU (Research Experiences for Undergraduates) and NREUP (National Research Experience for Undergraduates Program). Most universities in United States have dedicated offices of undergraduate research. Those offices, as well as many academic departments provide local funding opportunities.

The primary goal of these programs is to encourage undergraduate students to seek graduate level training through masters and doctoral programs, thereby producing a better equipped future work force. There are special opportunities for underrepresented groups (women and minorities), so that students from these groups can compete with those who were fortunate enough to have better access to corresponding resources.

This paper highlights research capabilities of undergraduate students with varying backgrounds. All authors of this paper are heavily involved in undergraduate research in addition to their own primary research programs. In Section 2, we describe five different research projects successfully completed by undergraduate students under our mentorship. In Section 3, we provide some concluding remarks.

2. Examples of Undergraduate Research Projects

2.1. Variations of the Greenberg Binary Unrelated Question RRT Model

This work was presented at the annual conference of the Society of Statistics, Computer and applications held at Jammu, India during March 6-8, 2017. The work was done by an undergraduate student, David Suarez, at UNC Greensboro, under the direction of Dr. Sat Gupta. The work is since accepted for publication and will appear in Suarez and Gupta (2018).

This study relates to the Randomized Response Models introduced originally by Warner (1965). The primary purpose of such models is to provide respondents in a survey enough privacy so that they can answer even sensitive and personal questions in a face-to-face survey. Warner's approach was to mix randomly direct and indirect versions of a sensitive question. For example, some of the respondents in the survey may face the direct question "Did you file an incorrect income tax return last year?" while others may face the indirect version of the question "Did you file a correct income tax return last year?". Clearly both questions reflect on the proportion of tax payers who cheat on taxes. The researcher does not know which question, direct or indirect, is being answered by the respondent. And there lies the anonymity.

Greenberg et al. (1969) offered an alternative model, ostensibly to offer greater anonymity to respondents. In Warner's model, all respondents answered a sensitive question, directly or indirectly. In Greenberg et al. (1969), some of the respondents, by random selection, answer an unrelated question. Let us recall that model first.

Let π_x be the unknown prevalence of a sensitive attribute X in the population and π_y be the prevalence of a non-sensitive attribute Y . A randomization device offers respondents a choice between two questions – a sensitive question with probability p (assumed known) and a non-sensitive unrelated question with probability $(1-p)$. Let p_y be the probability of a "yes" response. Then

$$p_y = \pi_x p + \pi_y (1 - p),$$

which leads to an estimator of π_x given by

$$\hat{\pi}_x = \frac{\hat{p}_y - \pi_y(1-p)}{p},$$

where \hat{p}_y is the sample proportion of “yes” responses. Variance of this estimator is given by

$$\text{Var}(\hat{\pi}_x) = \frac{p_y(1-p_y)}{np^2}.$$

In Suarez and Gupta (2018), we considered several variations of this model. In one of the variations, we allow a respondent to provide multiple independent responses. In another variation we use the inverse sampling technique, stopping at the k th “yes” response. It turns out that both of these techniques lead to more efficient estimates of the prevalence of the sensitive characteristic in the population of interest. The main result of the paper was

$$\begin{aligned} \text{Var}(\hat{\pi}_{G_k}) &< \text{Var}(\hat{\pi}_{G_1}) \text{ if } k > 1 \\ &< \text{Var}(\hat{\pi}_{GM}) \text{ if } mp_y < 1 \\ &< \text{Var}(\hat{\pi}_G) \text{ if } m > 1, \end{aligned}$$

where $\hat{\pi}_G$ is the usual Greenberg et al. (1969) estimator described above, $\hat{\pi}_{GM}$ is the estimator based on m independent responses, $\hat{\pi}_{G_1}$ is the estimator based on inverse sampling with stopping at the first “yes” response, and $\hat{\pi}_{G_k}$ is the estimator based on inverse sampling, stopping at the k th “yes” response. Also m is the number of independent responses used in $\hat{\pi}_{GM}$. Extensive simulation results validated the theoretical findings. The student researcher is now pursuing a master’s degree in statistics.

2.2. Predicting the Incidence of Sexually Transmitted Disease with Optional Greenberg Unrelated Question Model

Two other undergraduate students worked under the direction of Dr. Sat Gupta and introduced the optional version of the Greenberg et al. (1969, 1971) models for the binary and quantitative response situations. In the quantitative response model of Greenberg et al. (1971), a respondent answers the real question with probability p and an unrelated question with probability $(1-p)$. In Optional RRT models, introduced by Gupta et al. (2002), the main idea is that a question may be very sensitive for one person but not for another, so why force everyone to provide a scrambled response. In such models, a respondent is asked to give a truthful response if he/she does not consider the question sensitive, and a scrambled response otherwise, using Greenberg et al. (1969) model in the binary case, and the Greenberg et al. (1971) model in the quantitative case. The student researchers, Anna Tuck and Tracey Spears-Gill introduced the optional versions of the Greenberg et al. (1969, 1971) models in Gupta et al. (2013) and used these models to estimate the prevalence of STD among UNC Greensboro students in Spears-Gill et al. (2013).

The probability of a “yes” response in the binary optional model is given by

$$P_y = (1-w)\pi + w\{\pi p + (1-p)\pi_a\},$$

where π_a is the known prevalence of an unrelated characteristic, π is the unknown proportion of population that belongs to the sensitive group, p is the known probability of the respondent selecting the sensitive question, and w is the sensitivity level of the survey question in the population, that is, a proportion w of the respondents consider the question sensitive and hence choose to provide a scrambled response.

The corresponding quantitative response model is given by

$$Z = \begin{cases} X & \text{with probability } (1-w) + wp \\ Y & \text{with probability } w(1-p) \end{cases},$$

where Z is the reported response, X is the true response, Y is the response to the unrelated question, p is the known probability of the respondent selecting the sensitive question, and w is the sensitivity level of the survey question in the population.

It is easy to estimate π , μ_x and w from a random sample of responses from these models. The theoretical models from this work were later extended by Claudia Leonardi at LSU Health Sciences Center, New Orleans to earn a doctoral degree. Another student, Anu Chhabra, extended this work in a different direction to earn a PhD degree from Delhi University.

Anna Tuck is currently pursuing a PhD degree in Biostatistics at University of Washington-Seattle, and Tracey Spears-Gill earned a Master's degree in Biostatistics from UNC Chapel Hill.

2.3. The Invasion of Asian Carp of the Upper Mississippi River

This project was completed as part of the collaborative efforts between Jan Rychtář (UNCG) and Hyunju Oh (Bennett College). It was funded by the National Research Experience for Undergraduates Program (NREUP). The results were presented at the 2013 NREUP research symposium, at 2013 UNCG RMSC conference, at 2013 Joint REU Meeting at NC A&T, at the 2014 Interdisciplinary Research Day at Bennett College, and at the 2014 Joint Mathematics Meetings at Baltimore, MD. The work has been published in the Springer Proceedings in Mathematics and Statistics (Everett et al. 2015).

Asian Carps were imported from China in the 1970s to improve water quality of aquaculture and to control aquatic vegetation ponds. However, they subsequently migrated from ponds into the Mississippi river, where they quickly reached high population density. They are invasive species, highly detrimental to ecological balance as they threaten the native fish. Three math/stats majors selected this project. One was a rising sophomore and two were rising juniors. The group worked well together with the more mature students helping the younger one. Being math majors, students loved to do the literature search of many biological facts. Students learned how to translate real word problem into mathematics, how to solve the math problem and then how to translate the solution back into the real world. They captured the most important features by selecting only three types of fish (native predator – Largemouth Bass, native prey – Gizzard Shad, and invasive Silver Carp). The students were highly motivated and well prepared to build the mathematical model and constructed the following deterministic model for the net benefits (benefits – costs) to all species.

$$\begin{aligned} N_{SC} &= \frac{w_{SC}}{p_{SC}w_{SC} + p_{GS}w_{GS} + p_{LB}w_{LB}} R + E_{SC} - C_{SC}^H - p_{SC}C_{SC}^M \\ N_{LB} &= \frac{w_{LB}}{p_{SC}w_{SC} + p_{GS}w_{GS} + p_{LB}w_{LB}} R + E_{LB} + p_{GS}V_{LB}^{GS} \\ N_{GS} &= \frac{w_{GS}}{p_{SC}w_{SC} + p_{GS}w_{GS} + p_{LB}w_{LB}} R + E_{GS} - p_{LB}C_{GS}^{LB} \end{aligned}$$

A stochastic version of this model can be expressed by adding error components in the model.

Even in this simplified model, they had to work with 14 parameters. In the equations above, N_X is the net benefit to species X , w_X is the weight of species X , p_X is the prevalence of species X , E_X is the egg production of species X , R is the amount of resources in the river, C_{SC}^H is the cost of human measures to Silver Carp, C_{SC}^M is the cost of migration of Silver Carp, V_{LB}^{GS} is the value of caught Gizzard Shads to Largemouth Bass, and C_{GS}^{LB} is the predator cost Largemouth Bass causes to Gizzard Shad.

Students performed the stability analysis and studied conditions under which the Silver Carp cannot invade the native species. As one of the conclusions, we showed that the invasion can be best prevented by lowering the amount of available resources in the river while at the same time selectively harming the population of Silver Carps. Students also performed Monte Carlo simulations to make sure the results are robust.

This project brought some frustration to the students when we discovered an issue with the model after which a significant part had to be redone from the beginning. After initial discouragement, the students found it to be quite a valuable experience that distinguished this research from regular class work with only clear and straightforward answers.

2.4. k -kings in Products of Digraphs

This project in graph theory was completed by undergraduate students Morgan Morge, a senior applied mathematics major, and Peter LaBarr, a rising senior majoring in secondary mathematics teaching. The students worked in a multi-leveled research group, which included a doctoral student in operations research to help with the learning of introductory graph theoretic material, LaTeX, and general presentation skills, under the supervision of faculty Dewey Taylor at Virginia Commonwealth University. The research group met as a team for 3 hours each week during the entire academic year.

A directed graph (or digraph) consists of a non-empty finite set of elements called vertices and a finite set of ordered pairs of distinct vertices called arcs. Typical applications of digraphs include analyzing social networks, finding shortest routes, creating project schedules and analyzing network flows. Digraphs have also been used to represent casual or temporal relationships and came to be known as Bayesian networks (Pearl, 1985).

A k -king in a digraph D is a vertex that can reach every other vertex in D by a directed path of length at most k . The notion of a k -king is a generalization of a king in a tournament. There are four standard digraph products (Cartesian, direct, strong, and lexicographic), each having a long history with numerous applications and open problems (Hammack et al., 2011). The goal of this project was to classify k -kings in all four standard digraph products. In particular, we determined the relationship between k -kings in the product of digraphs and k -kings in the factors of the product. The four results we proved are summarized in the Theorem below.

Theorem Let D_1 and D_2 be digraphs.

1. (x_1, x_2) is a k -king in the strong product of D_1 and D_2 if and only if each x_i is a k -king in D_i .
2. (x_1, x_2) is a k -king in the Cartesian product of D_1 and D_2 if and only if each x_i is a k_i -king in D_i , where $k_1 + k_2 = k$.
3. (x_1, x_2) is a k -king in the lexicographic product of D_1 and D_2 if and only if x_1 is a k -king in D_1 , and either x_2 is a k -king in D_2 or the length of the shortest directed cycle in D_1 containing x_1 has length k or less.

4. If (x_1, x_2) is a k -king in the direct product of D_1 and D_2 , then each x_i is a k -king in D_i .

These results are currently being published in Hammack (in press). The two undergraduates presented their work at the annual UNCG Regional Mathematics and Statistics Conference in November 2016 and at the Virginia Academy of Sciences Conference in May 2017.

2.5. Game-theoretic Models of Individual Vaccination Decisions

Bauch and Earn (2004) developed a framework for applying game-theoretic models to vaccination decisions. This framework had been used to address voluntary vaccination policies for smallpox (Bauch et al. 2003), influenza (Galvani et al. 2007), measles (Shim et al. 2012), rubella (Shim et al. 2009), and toxoplasmosis (Sykes and Rychtář 2015). This general framework can also be adapted to other personal protection decisions, such as the use of insecticide-treated cattle to eliminate African sleeping sickness (Crawford et al. 2015), mosquito repellent to combat dengue fever (Dorsett et al. 2015), and insecticide-treated bed-nets to fight malaria (Broom et al. 2016).

In a vaccination game, an individual has two strategies: to vaccinate or to not vaccinate. Each strategy is assigned a corresponding expected payoff:

$$E_v = -C - \pi_v,$$

$$E_{nv} = -\pi_{nv},$$

where C is the cost of vaccination relative to the cost of infection, π_v is the probability of getting infected if vaccinated, and π_{nv} is the probability of getting infected if not vaccinated. The probability of getting infected depend on the force of the infection, which in turn depends on vaccination decisions of other individuals. The goal is to find a Nash equilibrium solution to this game: this is a population vaccination rate such that no individual can improve its payoff by deviating from this optimal strategy.

In summer 2016, three undergraduate student participants of the math-bio REU program at UNCG: Julia Kobe (Wentworth Institute of Technology), Neil Pritchard (UNCG), and Ziaqueria Short (Winston-Salem State University), investigated vaccination and clean water usage as optimal personal protection strategies against cholera. The project was jointly supervised by Igor Erovenko and Jan Rychtář. The students found that while it is impossible to reach herd immunity through voluntary policies, optimal personal protection levels ensure very low endemic levels of cholera provided the cost of protective measures is sufficiently low relative to the cost of getting infected.

Neil Pritchard presented these results at the Eighth Annual Undergraduate Research Conference at the Interface of Biology and Mathematics at NIMBioS, University of Tennessee, in October 2016. This research resulted in a paper that has been submitted for publication in a peer-reviewed journal. The paper is currently under review. Julia Kobe had been accepted to a summer research program at Harvard University in 2017.

In summer 2017, Igor Erovenko supervised two groups of undergraduate student participants of the math-bio REU program at UNCG:

1. Andrew Brettin (University of Minnesota), Rosa Rossi-Goldthorpe (Bowdoin College), and Kyle Weishaar (Regis University)
2. Jonathan Machado (UNCG), Alfredo Martinez (Whitworth University), and Eric Sanchez (UNCG)

These students chose to investigate optimal vaccination strategies for Ebola and meningitis, respectively. The first group established a remarkable result: eradication of Ebola

through voluntary vaccination could be possible if it is coupled with educational efforts regarding the effectiveness of the vaccine and dangers of getting infected with Ebola virus. Since the Ebola vaccine is in the last stages of testing and about to be released to the public, such a conclusion may have a profound influence on how it is accepted.

Both groups also performed uncertainty and sensitivity analysis of the outcomes on the choice of the model parameters. We used a Latin hypercube sampling method to effectively sample the parameter space. The sampled points were then used to compute the partial rank correlation coefficients (PRCC) to determine which parameters have the highest influence on the model outcomes.

Utilizing these statistical techniques becomes standard and expected practice for such problems. Parameters of epidemiological models are often chosen so as to fit real world data. However, such models often contain dozens of different parameters. Therefore, we can't be sure that each parameter was assigned an accurate value. In uncertainty and sensitivity analyses, parameters of the model vary over a certain range of values, rather than being assigned one fixed value. This allows to perform rigorous statistical analysis of possible outcomes if different values of parameters are chosen. In particular, the first group of students working on Ebola demonstrated that their conclusion -- Ebola could be eradicated through voluntary vaccination -- does not depend on the choice of the epidemiological model parameters. In this case, applying statistical techniques to the game-theoretic model outcomes significantly strengthened the plausibility of the result.

Both student groups prepared draft versions of papers based on their work. These papers should be submitted for publication in the near future.

3. Concluding Remarks

We have shown five different research projects, each one leading to conference presentation(s) and publication(s) in peer reviewed journals and books. The diversity of projects in Section 2 clearly demonstrates that, with proper mentoring, undergraduate students can be successful in many different areas of mathematics and statistics. Students can work on open research problems and make substantial progress towards solutions.

An important factor that allows students to focus on research is funding. Nevertheless, we believe that it is the faculty mentors, not the funding, that plays the most crucial role in students' success. Faculty mentors need to select appropriate research problems that are challenging enough to ignite the natural intellectual curiosity of the students, yet do not require lots of background knowledge. Problems that students can get started on quickly lead to initial successes that give students the confidence and interest to continue. In addition, faculty mentors need to believe that their students can be successful and achieve their goals. Mentors should be supportive and provide the right help at the right time.

Good mentorship comes with experience and faculty should be prepared to make a few mistakes along the way as well. We encourage all faculty to give undergraduate research a try. Mentoring undergraduate students is very rewarding and is often easier than one may think. You will be amazed how creative and hardworking some students can be outside the standard classroom setting.

Acknowledgements

Projects 2.2 – 2.5 were supported by various grants which we would like to acknowledge. NREUP is an MAA activity funded by a NSF grant DMS-1359016. REU program was funded by the NSF grants DMS-1359187 and DMS-1659646. The 4th author Rychtar will also like to acknowledge Simons Foundation grant 245400.

References

- Bauch, C.T. and Earn, D.J.D. (2004). Vaccination and the theory of games. *Proceedings of the National Academy of Sciences*, **101**,13391-13394.
- Bauch, C.T., Galvani, A.P. and Earn, D.J.D. (2003). Group-interest versus self-interest in smallpox vaccination policy. *Proceedings of the National Academy of Sciences*, **100**, 10564-10567.
- Broom, M., Rychtář, J. and Spears-Gill, (2016). The game-theoretical model of using insecticide treated bed-nets to fight malaria. *Applied Mathematics*, **7**, 852-860.
- Crawford, K., Lancaster, A., Oh, H. and Rychtář J. (2015). A voluntary use of insecticide-treated cattle can eliminate African sleeping sickness. *Letters in Biomathematics*, **2**(1), 91-101.
- Dorsett, C., Oh, H., Paulemond, M.L. and Rychtář, J. (2015). Optimal repellent usage to combat dengue fever. *Bulletin of Mathematical Biology*, **78**(5), 916-922.
- Everett, J., Jasim, M., Oh, H., Rychtář, J. and Smith, H. (2015). Modeling the Asian Carp Invasion Using Mathematical Evolutionary Game Theory, *Springer Proceedings in Mathematics and Statistics*, Vol. **109**, 81-90.
- Galvani, A.P., Reluga, T.C. and Chapman, G.B. (2007). Long-standing influenza vaccination policy is in accord with individual self-interest but not with the utilitarian optimum. *Proceedings of the National Academy of Sciences*, **104**, 5692-5697.
- Gargeya, M-S., and Pratap, P. (2015). Exploration into the Harmonic Structure of the Tabla, *Collaborative Mathematics and Statistics Research. Springer International Publishing*, 3-6.
- Goel, S. and Tashakkori, R. (2015). Correlation Between Body Measurements of Different Genders and Races. *Collaborative Mathematics and Statistics Research. Springer International Publishing*, 7-17.
- Greenberg, B.G., Abul-Ela, A.L.A., Simmons, W.R., and Horvitz, D.G. (1969). The unrelated question randomized response model: theoretical framework. *Journal of the American Statistical Association*, **64**, 520-529.
- Greenberg, B.G., Keublar, R.T., Abernathy, J.R., and Horvitz, D.G., (1971). Application of randomized response technique in obtaining quantitative data. *Journal of the American Statistical Association*, **66**, 243-250.
- Gupta, S., Gupta, B. and Singh, S. (2002). Estimation of the Sensitivity Level of Personal Interview Survey Questions. *Journal of Statistical Planning and Inference*, Vol. **100**, pp. 239-247.
- Gupta, S., Tuck, A., Spears Gill, T., and Crowe, M. (2013). Optional Unrelated-Question Randomized Response Models. *Involve: A Journal of Mathematics*, Vol. **6** (4), 483-492.
- Hammack, R. (in press). Digraph Products, in Bang-Jensen, Jørgen, and Gregory Z. Gutin (Eds). *Classes of Directed Graphs*, Springer-Verlag.
- Hammack, R., Imrich, W. and Klavžar, S. (2011). *Handbook of product graphs*. Second edition, CRC press, Boca Raton, FL.

- Lopatto, D. (2007). Undergraduate research experiences support science career decisions and active learning. *CBE life sciences education*. 2007;6(4):297-306. Epub 2007/12/07. doi: 10.1187/cbe.07-06-0039. PubMed PMID: 18056301; PubMed Central PMCID: PMC2104507.
- Lopatto, D. (2009). *Science in Solution: The Impact of Undergraduate Research on Student Learning*. Tucson, AZ: Research Corporation for Science Advancement.
- Pearl, J. (1985). Bayesian networks: A model of self-activated memory for evidential reasoning. *Proceedings of the 7th Conference of the Cognitive Science Society*.
- Russell, S.H., Hancock, M.P. and McCullough, J. (2007). The pipeline. Benefits of undergraduate research experiences. *Science* 316(5824):548-9. Epub 2007/04/28. doi: 10.1126/science.1140384. PubMed PMID: 17463273.
- Seymour, E., Hunter, A.B., Laursen, S.L., DeAntoni, T. (2004) Establishing the Benefits of Research Experiences for Undergraduates in the Sciences: First Findings from a Three-Year Study. *Science Education*. **88(4)**, 493-534. PubMed PMID: EJ759873.
- Shim, E., Grefenstette, J.J., Albert, S.M., Cakourus, B.E. and Burke, D.S. (2012). A game dynamic model for vaccine skeptics and vaccine believers: measles as an example. *Journal of Theoretical Biology*, **295**, 194-203.
- Shim, E., Kochin, B. and Galvani, A. (2009). Insights from epidemiological game theory into gender-specific vaccination against rubella. *Mathematical Biosciences and Engineering*, **6(4)**, 839–854.
- Spears-Gill, T., Tuck, A., Gupta, S., Crowe, M. and Figueroa, J. (2013). A Field Test of Optional Unrelated Question Randomized Response Models: Estimates of Risky Sexual Behaviors. *Topics from the 8th Annual UNCG Regional Mathematics, and Statistics Conference, Springer Proceedings in Mathematics and Statistics Series*, Vol. **64**, 135-146.
- Suarez, D and Gupta, S. (2018). Variations of the Greenberg Unrelated Question Binary Model. *Involve: A Journal of Mathematics*, Vol. **11(1)**, 119-126. DOI: 10.2140/involve.2018.11.119
- Sykes, D. and Rychtář, J. (2015). A game-theoretic approach to valuating toxoplasmosis vaccination strategies. *Theoretical Population Biology*, **105**, 33-38.
- Warner, S.L. (1965). Randomized response: A survey technique for eliminating evasive answer bias. *Journal of American Statistical Association*, **60(309)**, 63-69.