

Is There a Cost to Convenience? An Experimental Comparison of Data Quality in Laboratory and Online Studies

Scott Clifford* and Jennifer Jerit†

Abstract

Increasingly, experimental research is being conducted on the Internet in addition to the laboratory. Online experiments are more convenient for subjects and researchers, but we know little about how the choice of study location affects data quality. To investigate whether respondent behavior differs across study location, we randomly assign subjects to participate in a study in a laboratory or in an online setting. Contrary to our expectations, we find few differences between participants in terms of the level of attention and socially desirable responding. However, we find significant differences in two areas: the degree of self-reported distractions while completing the questionnaire and the tendency to consult outside sources for answers to political knowledge questions. We conclude that when the greater convenience (and higher response rates) of online experiments outweighs these disadvantages, Internet administration of randomized experiments represent an alternative to laboratory administration.

Keywords: Mode effects, political knowledge, lab experiment, online experiment, attention, social desirability.

In their landmark textbook on experimental political science, Morton and Williams describe *location* as “probably the most salient dimension over which experiments differ” (2010, 278). Recent research finds an “inherent advantage” for questionnaire self-administration on the computer (Chang and Krosnick 2010), but computer administration can take place in a variety of places, such as a university laboratory or a respondent’s home. Indeed, the choice of where to conduct a study is a crucial consideration for researchers conducting an individual decision-making experiment, with a growing number of scholars turning to the Internet (Sargis et al. 2014). In this

The authors thank Kevin Arceneaux, Ryan Enos, Yanna Krupnikov, Thomas Leeper, Kerri Milita, Spencer Piston, John Barry Ryan, and the anonymous reviewers for helpful comments and suggestions. They also thank Jason Barabas for assisting with the design of the political knowledge questions.

*Department of Political Science, University of Houston, Houston, TX, USA;
e-mail: scott.clifford@duke.edu; scottaclifford@gmail.com

†Department of Political Science, Stony Brook University, Stony Brook, NY, USA;
email: jennifer.jerit@stonybrook.edu

study, we examine respondent behavior in a laboratory versus an online context—one of the first such mode comparisons of which we are aware.¹

The choice of lab versus online administration has important consequences for the burden on subjects and researchers, as well as the quality of the data. From a practical standpoint, online questionnaires can be distributed easily via email, Web sites, or crowd-sourcing platforms such as Mechanical Turk, and participants can complete the survey at a time and place of their choosing (Berinsky et al. 2012; Cassese et al. 2013). By contrast, laboratory experiments involve greater costs in terms of the administration of the study (e.g., setup, proctoring) and the potential inconvenience for subjects taking the study at a specified time and place. If there were few differences between data obtained in the lab versus the Internet (e.g., in terms of the quality of responses), researchers might conduct more of their experiments, even those involving college students, through an online platform. From a theoretical standpoint, this study extends research on the generalizability of findings across experimental contexts (e.g., Coppock and Green 2013; Jerit et al. 2013). Lab and online experiments vary in ways that affect whether the treatment is received. Thus, the decision to administer an experiment online or in the laboratory may have implications for the conclusions one draws from such studies.

EXPECTATIONS

Drawing upon Jerit et al. (2013), we expect that experiments administered in a laboratory and online setting will differ principally in terms of experimenter control and the obtrusiveness of the setting. In a lab, subjects complete the study under the discretion of the researcher and at a common location. In an online setting, subject interaction with the researcher is indirect. There also is more “behavioral latitude” (Gerber 2011, 120) in terms of what a subject does while completing a questionnaire online and greater noise from the outside world. As a result of these differences, we expect that participants in a lab study will devote higher levels of attention to the task than online participants (Hypothesis 1). Previous research shows that the mode of administration is related to social desirability pressures (Tourangeau and Yan 2007). For example, questionnaires delivered over the Internet (i.e., self-administration) exhibit lower levels of socially desirable responding than those delivered over the telephone (i.e., aural administration; Chang and Krosnick 2009). But self-administration may take place in the lab or online, making the implications for socially desirable responding somewhat unclear. Following the logic of Chang and Krosnick (2009; also see Evans et al. 2003), we surmise that participants in a lab may be more concerned with impression management than online participants because the trappings of a scientific study are more apparent (e.g., a proctor, other participants). This leads to the expectation that the level of socially

¹At the time of this writing, only Evans et al. (2003; Study 2) and Weigold et al. (2013; Study 1) had compared respondents in lab and online settings.

desirable responding will be higher in the lab compared with the online setting (Hypothesis 2).²

As more researchers use the Internet to collect data, there is growing concern that the behavioral latitude in online studies leads subjects to cheat on knowledge questions by consulting the Internet for answers (e.g., Vavreck 2012; Warren 2012). Others, citing the tendency of respondents to satisfice, doubt that people are sufficiently motivated to cheat. Thus, we also collected data on respondents' political knowledge. Our study is uniquely situated to investigate this issue because participants were sampled from the same target population and then randomized to either a lab or online condition after agreeing to participate in the study. Thus, any differences in the observed levels of political knowledge across the two conditions can be attributed to features of the experimental setting.

EXPERIMENTAL DESIGN

Respondents

Participants ($n = 435$) in the study were undergraduate students enrolled in political science classes at a large public university in the south in the spring of 2013.³ They were recruited to participate in exchange for extra credit and instructed to sign up for an appointment through a link hosted at the Department of Political Science Web site. All subjects, irrespective of their eventual treatment assignment, signed up for the study through the same mechanism (believing the study would take place in a computer lab on campus). The participants included 201 males and 234 females, with nearly half indicating they were either in their first or second year of school. Approximately 72% of the sample was White; 8% were African American, and 13% were Hispanic.

Procedure

Our study is a between-subjects design with two conditions, lab versus online administration. Participants were randomized into condition by using the list of students who had signed up for an appointment in combination with a random number generator. Depending on treatment assignment, participants were instructed (via e-mail) to come to a computer lab at a particular time during a five-day period, or they were told they would be receiving a link to a survey that they could complete at a time and place of their choosing during this same five-day period. Because participants were randomly assigned to condition after they had

²Tourangeau and Yan (2007, 869–870) note that the physical presence of the interviewer may not matter as much as the perception that one's responses are anonymous (but see Lelkes et al. 2012, on the effects of anonymity).

³This study was approved by the Human Subjects Review Committee at Florida State University (Application no. HSC 2013-10185).

already signed up for the study, any observed differences between responses in the lab and online conditions are likely due to the effects of the experimental context, rather than the differences between participants in each setting. Table A1 shows that demographic and other characteristics were similar across experimental conditions.

Measures

We investigated whether there were differences across mode of administration in three areas relating to data quality: respondent attention, socially desirable responding (SDR), and levels of political knowledge. Overall, seven questions were used to measure respondent attention: an instructional manipulation check or IMC (Oppenheimer et al. 2009), two bogus items (e.g., Meade and Craig 2012), two substantive manipulation checks that followed experimental treatments appearing elsewhere on the questionnaire, and two self-report items.⁴ Whereas IMCs are considered a general measure of attention (Berinsky et al. 2014), substantive manipulation checks determine whether a particular experimental treatment was received (Mutz 2011). The self-report measures asked individuals to assess their level of attention during the study. The first item asked respondents to indicate how closely they were paying attention to the questions (e.g., Berry et al. 1992). The second asked them to indicate which of several different activities they engaged in while answering the questionnaire (e.g., using a cell phone). Socially desirable responding is measured with a three-item battery from the Modern Racism Scale (McConahay 1986) and an open-ended item asking respondents how many sexual partners they have had in their lifetime. Finally, we measure political knowledge with eight questions about current events.⁵

EMPIRICAL RESULTS

Response Rates

We begin by describing the response rates across conditions. There was a significantly higher response rate in the online condition compared with the lab condition (88% vs. 77%, $p < .001$), which may reflect the greater convenience of online administration.⁶ To rule out concerns about selection bias related to the higher response rate in the online condition, we examined whether online subjects were different in terms of demographic and attitudinal characteristics. Across a range of variables that may be related to differential participation (race, GPA, political interest, voter registration status, year in school) there were no significant differences

⁴Our outcomes were part of a multi-investigator questionnaire consisting of approximately 80 questions, many of which were unrelated to the present study. The Appendix provides the wording for the outcomes examined in this study.

⁵Respondents were given unlimited time to answer the knowledge items.

⁶Response rates were calculated as the percentage of respondents completing the survey that originally signed up to participate.

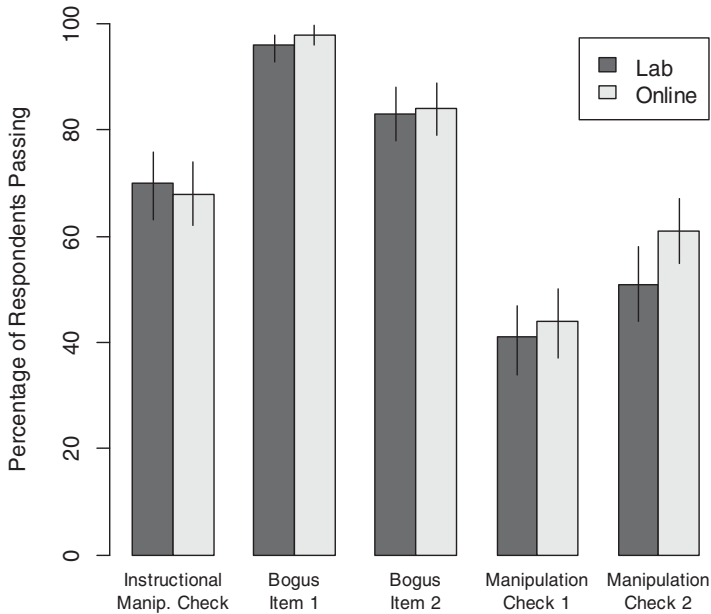


Figure 1

Respondent Attention by Mode of Administration

Note: Columns represent the percentage of respondents selecting the “correct” response to each question. Bars represent 95% confidence intervals.

between conditions (see Table A1 in the Appendix).⁷ Overall, there is little reason to suspect that the differential response rate led to selection effects across conditions. As a precaution, we confirmed that all of our results obtain in models including a basic set of controls.

Attention

Next we examine whether mode of administration affects attention to experimental stimuli and thus the likelihood of receiving the assigned treatment. The results of the five attention checks are summarized in Figure 1 (the self-report items are described separately). Starting on the left, passage rates for the instructional manipulation check are indistinguishable between the lab and online conditions (70% vs. 68%; $p = .75$). For the first bogus item, passage rates are above 95% in both conditions and thus are not significantly different from one another (96% vs. 98%; $p = .16$). For the second bogus item, there are no differences between the lab and online conditions (83% vs. 84%; $p = .82$). The first substantive manipulation check followed an experimental vignette in which respondents were randomly assigned to learn different facts about a politician’s earmarking activity. All subjects were told the politician had received earmarks and then, after two outcome measures, were asked

⁷There was one exception: Online subjects reported using the bus less often ($p = .08$).

Table 1
Scale Reliability by Experimental Condition

Scale	Lab	Online	Items
Modern Racism Scale	.83	.79	3
Political interest - time 1	.85	.86	3
Political interest - time 2	.85	.88	3
Political knowledge	.60	.61	8
Character trait evaluations	.86	.88	10
Need to evaluate	.68	.42	3
Average	.78	.74	

Note: Cell entries display Cronbach's Alpha coefficient for each scale by experimental condition. Political interest was measure at the beginning (time 1) and end of the survey (time 2). The third column lists the number of survey questions used to construct the scale.

whether the politician was part of a group of congressmen who had foregone earmarks for the past two years.

Passage rates for both groups were substantially lower than the other attention items, but no differences emerged between the lab and online conditions (41% vs. 44%; $p = .48$). The final substantive manipulation check followed an experimental vignette that randomized a politician's partisanship, issue stance, and explanation for the issue stance. After 13 outcome measures, respondents were asked to recall the politician's partisanship. Passage rates on the second manipulation check were slightly lower among students in the lab than in the online condition (51% vs. 61%; $p < .05$). On the basis of the analysis of the attention checks, there is no consistent effect of mode on attention among student participants taking the questionnaire in the lab versus online.

Previous researchers have used scale reliability as an indicator of respondent attention, finding higher reliabilities among more attentive respondents (Berinsky et al. 2014; Huang et al. 2012). Table 1 shows the Cronbach's alpha coefficients for six multi-item scales in the survey. The average scale reliability in the lab condition is .78 compared with .74 in the online condition, suggesting no difference between modes in terms of data quality.⁸

Although students in the online and lab conditions paid similar levels of attention to the stimuli, online participants might have faced more interruptions depending on where and when they completed the questionnaire. This appears to be the case according to Figure 2, which shows the rates of self-reported distraction across groups. Starting on the left, there is a significantly higher rate of distraction among online participants from cell phone use (21% vs. 9%; $p < .001$), surfing the internet (11% vs. 1%; $p < .001$), and talking with another person (21% vs. 2%; $p < .001$).

⁸Straight-lining can increase scale reliability even though it is a form of satisficing (Huang et al. 2012). We found no significant differences in straight-lining across conditions, and our scale reliability results are substantively unchanged by removing all straight-liners prior to analysis.

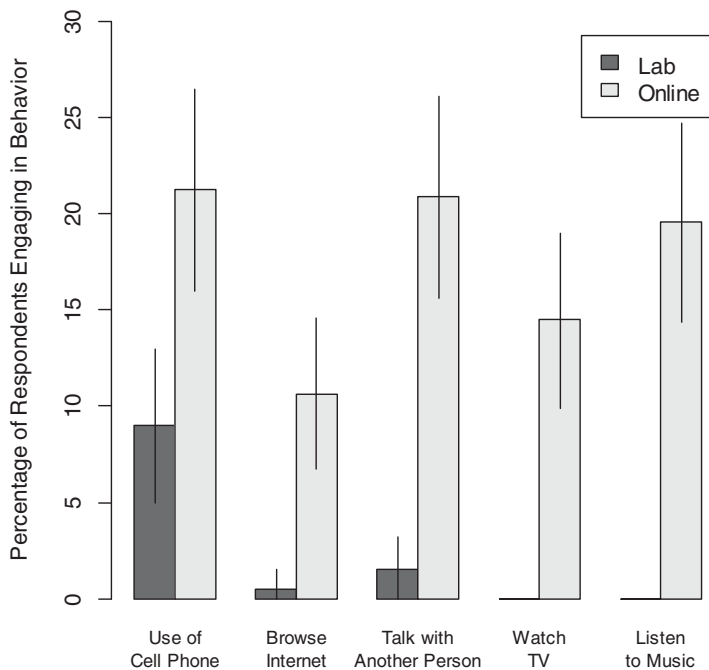


Figure 2

Respondent Distractions by Mode of Administration

Note: Columns represent the percentage of respondents engaging in each behavior. Bars represent 95% confidence intervals.

Online subjects also were asked about watching TV and listening to music. Lab subjects were not asked these questions because it would have been impossible for them to engage in these behaviors in the lab. Of online students, 14% reported watching TV during the survey and 20% said they were listening to music during the survey. Finally, in response to the item about paying attention to what the survey questions mean, students in the lab condition reported paying greater attention than students taking the questionnaire online (4.0 vs. 3.8; $p < .05$). Taken together, this evidence suggests that subjects in the online condition faced significantly higher rates of distraction from a number of sources (though these distractions were not associated with worse performance on the attention checks).

Socially Desirable Responding

We now examine the presence of socially desirable responding. As a first indicator of socially desirable responding, we analyze item non-response on two sensitive questions (Berinsky 1999; Tourangeau and Yan 2007). Contrary to our expectations, there was no missingness on the Modern Racism scale in either condition.⁹

⁹ Respondents were not forced to provide a response, but were notified if they had missed a question before continuing to the next page (a “requested response” in Qualtrics).

Additionally, there was a surprisingly low level of missingness on the sexual partners question across both conditions (lab: 3%; online: 2%), with no significant difference between groups ($p = .56$). Finding little evidence of non-response, we turn to subjects' responses to determine if mode affected self-reported opinions. There is no difference between conditions on the Modern Racism scale (4.1 vs. 4.2; $p = .58$), or the number of sexual partners (6.7 vs. 5.7, $p = .38$). Tourangeau and Smith (1996) report that although social pressure decreases the self-reported number of partners among women, it may increase self-reports among men. We investigated the mode effect separately by gender, but found no significant differences.¹⁰ Overall, there is little evidence that mode (lab vs. online) affects levels of socially desirable responding.¹¹

Knowledge

Our final analysis investigates levels of political knowledge across conditions. Figure 3 shows the distribution of correct responses (out of eight) for student participants in the lab and online conditions. Consistent with suspicions about cheating behavior in online surveys, students in the online condition scored significantly higher on the knowledge scale than lab students (6.4 vs. 5.9; $p < .01$). Indeed, 61% answered 7 or 8 questions correctly online, whereas only 44% of lab participants obtained a similar score. To buttress our claim that this difference stems from cheating, we examine the criterion validity of the knowledge scale by looking at its correlation with political interest. If online subjects are cheating, the knowledge scale should have lower criterion validity, as indicated by a weaker correlation with political interest (e.g., Prior 2009).¹² In line with the cheating interpretation, the correlation between interest and knowledge was higher in the lab ($r = .48$) than in the online condition ($r = .33$), a difference that is statistically significant ($p < .10$).¹³ We conclude that subjects in the online condition were more likely to cheat on knowledge items, weakening the validity of the scale. This interpretation also is consistent with Figure 2, which reveals that subjects in the online condition were more likely to report surfing the Internet than lab subjects (11% vs. 1%).¹⁴

¹⁰Men and women both reported fewer partners in the online condition, but neither of these differences were statistically significant (men: lab = 9.7, online = 8.4; women: lab = 3.9, online = 3.5).

¹¹This pattern contrasts with Evans et al. (2003), who find greater socially desirable responding in the lab versus online. In that study, the differences across settings was starker (e.g., the lab proctor wore a white lab coat, carried a clipboard, and sat in close proximity to subjects for the duration of the study).

¹²Recall from Table 1 that the political interest index was reliable across modes. In addition, mean values on the scale did not differ across modes (lab = 3.49, online = 3.54, $p = .50$).

¹³These results also hold with controls for demographic variables.

¹⁴We did not include reaction timers on the knowledge questions because their interpretation is ambiguous. Longer reaction times in the online condition might reflect the presence of more distractions, rather than cheating per se.

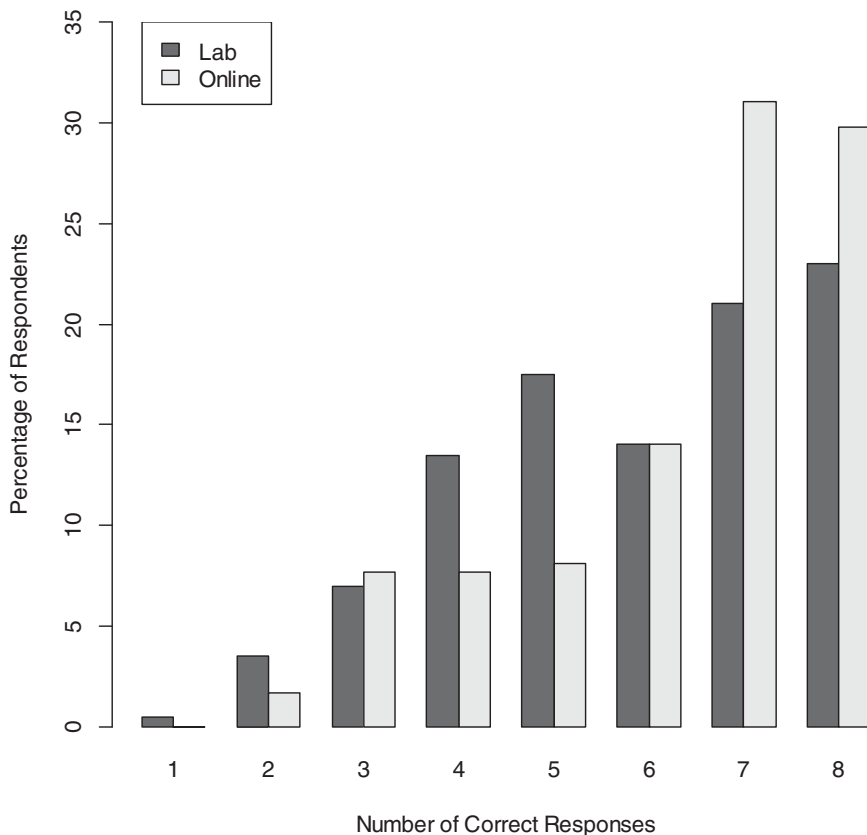


Figure 3

Political Knowledge Scores by Mode of Administration

Note: Columns represent the percentage of respondents in a particular category.

DISCUSSION

In a comparison of student subjects randomly assigned to take a questionnaire in the lab versus online, there were few differences in respondent attention across a variety of measures (contrary to our first hypothesis). From a practical standpoint, the results suggest that online experiments may be an appealing alternative to lab experiments. Online administration lowers the burden on the researcher and it appears to have a similar effect on subjects, as evidenced by the higher response rate in the online condition. That said, our results revealed substantially higher levels of distraction outside of the lab. This pattern is consistent with the idea that researchers lose control over key aspects of an experiment when a study takes place outside of the laboratory (McDermott 2002; Morton and Williams 2010). In our case, these distractions did not translate into worse performance on the attention checks, but our findings should give pause to those carrying out subtle or short-lived manipulations.

To further illustrate some of these challenges, consider the use of non-conscious primes, which are common in psychology and some subfields of political science (e.g., Bargh and Chartrand 2000; Lodge and Taber 2013). The brevity of the presentation ensures that the stimulus cannot be consciously processed, but potential distractions from an online setting might prevent the treatment from being received (though see Weinberger and Westen 2008, for an exception). In other instances, a concept or trait might be successfully primed in an online study (say, though a scrambled-sentence task), but its effects might not be observed if the subject becomes distracted by unrelated stimuli before answering the outcome measures.

Regarding our second hypothesis, mode does not appear to affect socially desirable responding. Across two topics previously shown to create social desirability pressures, we found no differences between experimental conditions, either in terms of the patterns of non-response or substantive responses. Previous research has shown that online surveys create weaker social desirability pressures relative to phone or face-to-face interviews. Self-administered questionnaires in a laboratory environment fare no worse on this dimension.

Finally, our results have important implications for research on political knowledge. We found evidence that students in our online condition were more likely to cheat on the knowledge items, generating higher knowledge scores and weakening the validity of the measure. It is unclear whether this pattern generalizes to other populations—non-students may not be as motivated to cheat on knowledge questions. Given the rise of online surveys, however, researchers may consider designing questionnaires in a way that discourages Internet surfing, particularly if political knowledge is the outcome of interest.¹⁵

Our study highlights some of the issues related to the cost and convenience of online versus lab studies, but cost and convenience are not the only considerations when conducting an experiment. Online administration can be advantageous when a researcher wants to collect data from a nonlocal sample (either national or international). Conversely, some studies are difficult, if not impossible, to administer online, such as those that collect physiological data or that involve human confederates as part of the treatment. Nevertheless, the growing use of the Internet as a platform for data collection points to a need for studies that explore mode differences between experiments conducted online and in the lab.

SUPPLEMENTARY MATERIAL

To view supplementary material for this paper, please visit <http://dx.doi.org/10.1017/xps.2014.5>.

¹⁵Some researchers include warnings that advise participants not to look up answers or get assistance from others (e.g., Berinsky et al. 2012; Boster and Shulman 2014; Goodman et al. 2013). Additionally, some survey software programs “take over” the respondent’s screen, making it difficult to open a new tab or browser to search for answers.

REFERENCES

- Bargh, J., and Chartrand, T. L. 2000. The Mind in the Middle: A Practical Guide to Priming and Automaticity Research. In *Handbook of Research Methods in Social and Personality Psychology* eds. H. T. Reis and C. M. Judd. New York: Cambridge University Press, 253–85.
- Berinsky, A. J. 1999. Can We Talk? Self-Presentation and the Survey Response. *Political Psychology* 25 (4): 643–59.
- Berinsky, A. J., Huber, G. A., and Lenz, G. S. 2012. Using Mechanical Turk as a Subject Recruitment Tool for Experimental Research. *Political Analysis* 20 (3): 351–68.
- Berinsky, A. J., Margolis, M., and Sances, M. 2014. Separating the Shirkers from the Workers? Making Sure Respondents Pay Attention on Internet Surveys. *American Journal of Political Science* Forthcoming.
- Berry, D. T. R., Wetter, M. W., Baer, R. A., Larsen, L., Clark, C., and Monroe, K. 1992. MMPI-2 Random Responding Indices: Validation Using a Self-Report Methodology. *Psychological Assessment* 4 (3): 340–45.
- Boster, F., and Shulman, H. 2014. *Political Knowledge Test Performance as a Function of Venue, Time Pressure, and Performance Norms*. Working Paper, North Central College.
- Cassese, E. C., Huddy, L., Hartman, T. K., Mason, L., and Weber, C. R. 2013. Socially Mediated Internet Surveys: Recruiting Participants for Online Experiments. *PS: Political Science and Politics* 46 (4): 775–84.
- Chang, L., and Krosnick, J. A. 2009. National Surveys via RDD Telephone Interviewing versus the Internet. *Public Opinion Quarterly* 73 (4): 641–78.
- Chang, L., and Krosnick, J. A. 2010. Comparing Oral Interviewing with Self-Administered Computerized Questionnaires. *Public Opinion Quarterly* 74 (1): 154–67.
- Coppock, A., and Green, D. P. 2013. *Assessing the Correspondence Between Experimental Results Obtained in the Lab and Field: A Review of Recent Social Science Research*. Working Paper, Columbia University.
- Evans, D. C., Garcia, D. J., Garcia, D. M., and Baron, R. S. 2003. In the Privacy of Their Own Homes: Using the Internet to Assess Racial Bias. *Personality and Social Psychology Bulletin* 29 (2): 273–84.
- Gerber, A. 2011. Field Experiments in Political Science. In *Handbook of Experimental Political Science* eds. J. Druckman, D. P. Green, J. H. Kuklinski, and A. Lupia. New York: Cambridge University Press, 115–38.
- Goodman, J. K., Cryder, C. E., and Cheema, A. 2013. Data Collection in a Flat World: The Strengths and Weaknesses of Mechanical Turk Samples. *Journal of Behavioral Decision Making* 26 (3): 213–24.
- Huang, J. L., Curran, P. G., Keeney, J., Poposki, E. M., and DeShon, R. P. 2012. Detecting and Detering Insufficient Effort Responding to Surveys. *Journal of Business Psychology* 27: 99–114.
- Jerit, J., Barabas, J., and Clifford, S. 2013. Comparing Contemporaneous Laboratory and Field Experiments on Media Effects. *Public Opinion Quarterly* 77 (1): 256–82.
- Lelkes, Y., Krosnick, J. A., Marx, D. M., Judd, C. N., and Park, B. 2012. Complete Anonymity Compromises the Accuracy of Self-Reports. *Journal of Experimental Social Psychology* 48: 1291–99.
- Lodge, M., and Taber, C. S. 2013. *The Rationalizing Voter*. New York: Cambridge University Press.

- Meade, A. W., and Craig, S. B. 2012. Identifying Careless Responses in Survey Data. *Psychological Methods* 17 (3): 437–55.
- McConahay, J. G. 1986. Modern Racism, Ambivalence, and the Modern Racism Scale. In *Prejudice, Discrimination, and Racism* eds. J. F. Dovidio and S. L. Gaertner. New York: Academic Press, 91–125.
- McDermott, R. 2002. Experimental Methods in Political Science. *Annual Review of Political Science* 5: 31–61.
- Morton, R. B., and Williams, K. C. 2010. *Experimental Political Science and the Study of Causality*. New York: Cambridge University Press.
- Mutz, D. 2011. *Population-Based Survey Experiments*. Princeton, NJ: Princeton University Press.
- Oppenheimer, D., Meyvis, T., and Davidenko, N. 2009. Instructional Manipulation Checks: Detecting Satisficing to Increase Statistical Power. *Journal of Experimental Social Psychology* 45: 867–72.
- Prior, M. 2009. Improving Media Effects Research Through Better Measurement of News Exposure. *Journal of Politics* 71 (3): 893–908.
- Sargis, E. G., Skitka, L. J., and McKeever, W. 2014. The Internet as Psychological Laboratory Revisited: Best Practices, Challenges, and Solutions. In *The Social Net: The Social Psychology of the Internet* ed. Y. Amichai-Hamberger. Oxford: Oxford University Press. In Press.
- Tourangeau, R., and Smith, T. W. 1996. Asking Sensitive Questions: The Impact of Data Collection Mode, Question Format, and Question Context. *Public Opinion Quarterly* 60 (2): 275–304.
- Tourangeau, R., and Yan, T. 2007. Sensitive Survey Questions. *Psychological Bulletin* 133 (5): 859–83.
- Vavreck, L. 2012. The Myth of Cheating on Self-Completed Surveys. <http://today.yougov.com/news/2012/04/17/myth-cheating-self-completed-surveys/> (Accessed June 10, 2013).
- Warren, J. 2012. Fake Orgasms and the Tea Party: Just Another Political Science Convention. *The Atlantic Online*. <http://www.theatlantic.com/politics/archive/2012/04/fake-orgasms-and-the-tea-party-just-another-political-science-convention/255909/> (Accessed June 10, 2013).
- Weigold, A., Weigold, I. K., and Russell, E. J. 2013. Examination of the Equivalence of Self-Report Survey-Based Paper-and-Pencil and Internet Data Collection Methods. *Psychological Methods* 18 (1): 53–70.
- Weinberger, J., and Westen, D. 2008. RATS, We Should Have Used Clinton: Subliminal Priming in Political Campaigns. *Political Psychology* 29 (5): 631–51.