

SHORT REPORT  

How Dropping Subjects Who Failed Manipulation Checks Can Bias Your Results: An Illustrative Case

Simon Varaine 



PACTE, Université Grenoble Alpes, CNRS, Sciences Po Grenoble (School of Political Studies), 38000 Grenoble, France
Email: simon.varaine@sciencespo-grenoble.fr

Abstract

Manipulations checks are postexperimental measures widely used to verify that subjects understood the treatment. Some researchers drop subjects who failed manipulation checks in order to limit the analyses to attentive subjects. This short report offers a novel illustration on how this practice may bias experimental results: in the present case, through confirming a hypothesis that is likely false. In a survey experiment, subjects were primed with a fictional news story depicting an economic decline versus prosperity. Subjects were then asked whether the news story depicted an economic decline or prosperity. Results indicate that responses to this manipulation check captured subjects' preexisting beliefs about the economic situation. As a consequence, dropping subjects who failed the manipulation check mixes the effects of preexisting and induced beliefs, increasing the risk of false positive findings. Researchers should avoid dropping subjects based on posttreatment measures and rely on pretreatment measures of attentiveness.

Keywords: manipulation checks; randomized experiments; survey experiments; causal inference; Type I error

Manipulations checks are postexperimental measures aiming at “ensuring that an experiment actually has been conducted (i.e., that the Independent Variable has been effectively manipulated)” (Sansone et al. 2003, p. 244). They typically take the form of comprehension questions immediately following the experimental treatment to check that subjects paid attention and understood the treatment (Kane and Barabas 2019). The inclusion of manipulation checks enters standards of best practices in experimental political science (Mutz and Pemantle 2015). They are particularly important to avoid Type II error – i.e. the false negative conclusion that the research hypothesis is wrong while it is actually true – in case of null results.

  This article has earned badges for transparent research practices: Open Data and Open Materials. For details see the [Data Availability Statement](#).

© The Author(s), 2022. Published by Cambridge University Press on behalf of American Political Science Association

A widespread practice is to exclude participants who failed manipulation checks in order to limit the analyses to subjects who understood the experimental prompt (Aronow et al. 2019). However, some studies have warned that this may bias the analyses (Aronow et al. 2019; Berinsky et al. 2014; see also Montgomery et al. 2018). Berinsky et al. (2014) showed that individual responses to screeners correlate with a range of personal characteristics. Dropping inattentive subjects may distort the sample to certain “races, ages, and levels of education”. More problematically, Aronow et al. (2019) demonstrated that this may bias the estimation of causal effects by creating an asymmetry between experimental arms.

This study offers a new illustration on how dropping subjects who failed manipulation checks may bias experimental results. Aronow et al. (2019) presented an illustrative experiment in which dropping subjects lead to underestimating the effect size of the treatment. The present study presents another experimental case in which dropping subjects increases the risk of Type I error when testing a hypothesis of interest – i.e. drawing a false positive conclusion that confirms the research hypothesis while it is actually wrong.

The experiment was conducted in an online survey filled during April 2019 by nationally representative samples from Denmark, France, Germany, Italy, Spain, and the Netherlands. A total of 3949 subjects participated in the experiment, based on the economic threat manipulation from Stenner (2005). Subjects were randomly assigned to one of two short fictional news stories about the national economic context, respectively, depicting an improving situation (*prosperity*) or a worsening situation (*decline*)¹. The initial purpose of the experiment was to test whether subjects express more nostalgia after the *decline* treatment compared to the *prosperity* treatment.

Following the treatment, subjects answered a manipulation check question: “According to the news story, the national economic situation is: Worsening/Stable/Improving”. Only 63% of subjects provided the correct answer regarding their experimental treatment – i.e. “Improving” in the *prosperity* treatment and “Worsening” in the *decline* treatment. More problematically, subjects who failed did apparently not respond at random, as shown in Figure 1. 30% of subjects declared that, according to the news story, the economic situation was improving, 27% that it was stable and 43% that it was worsening.

Prior to the experiment, subjects answered a question about their own perception of the economic situation: “Would you say that the economic situation now is better or worse to how it was 5 years ago?”. Subjects responded with a 11-point scale from 0 (“Much worse”) to 10 (“Much better”). Results from a two-way ANOVA indicate that responses to this question are significantly related to responses to the manipulation check, $F(2,3731) = 55.79, p < .0001$. As shown in Figure 2, the more favorably subjects perceived the economy, the less they responded that the news story depicted an economic decline and the more they responded that the story depicted economic stability². This means that the manipulation check actually captured some subjects’ preexisting beliefs about the economic situation.

¹See the contents of treatments in the online appendix.

²In contrast, there is no clear effect on the probability that subjects responded that the treatment described a economic prosperity.

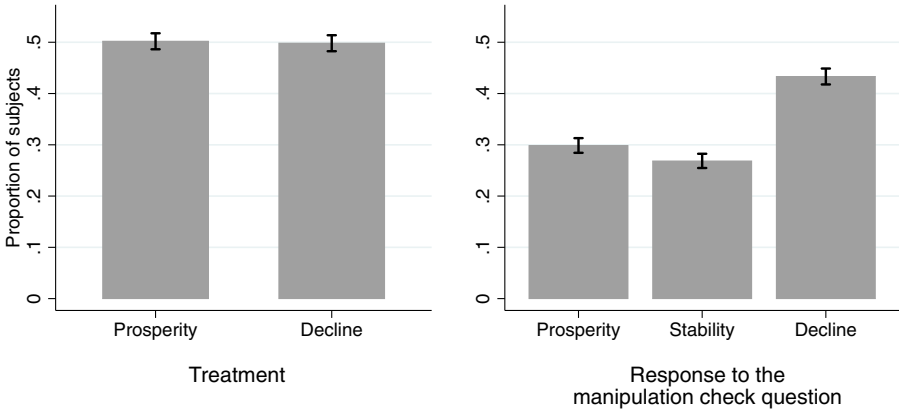


Figure 1

Share of subjects by experimental treatment and by response to the manipulation check question (with 95% Confidence Interval).

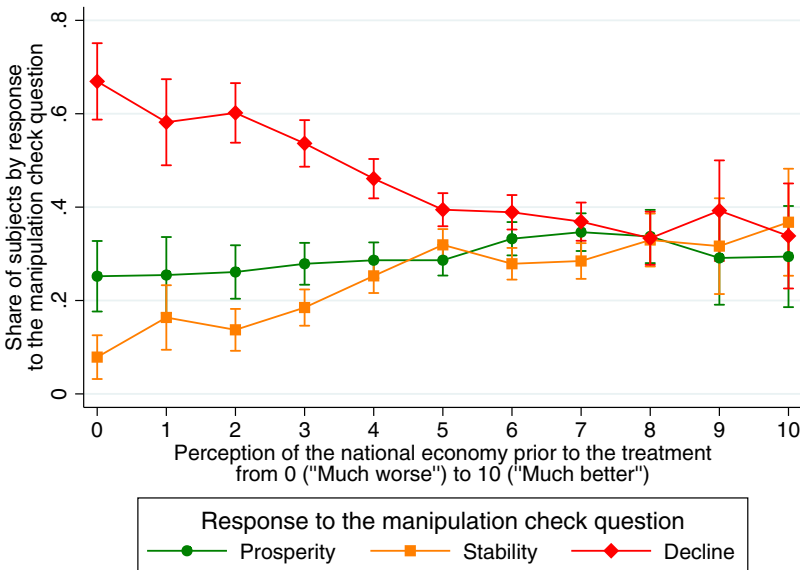


Figure 2

Share of responses to the manipulation check question depending on the perception of the national economy prior to the experiment (with 95% Confidence Interval).

Now, what happens if we drop subjects who failed the manipulation check? Figure 3 presents the average perception of the economy prior to the experiment for subjects of the *prosperity* versus *decline* treatments. When including all subjects, there is no significant difference across treatments in the perception of the national economy prior to the survey experiment, $t(3752) = 0.6062, p = .5444$. This is what we expect from randomization: the treatment is independent from the subjects

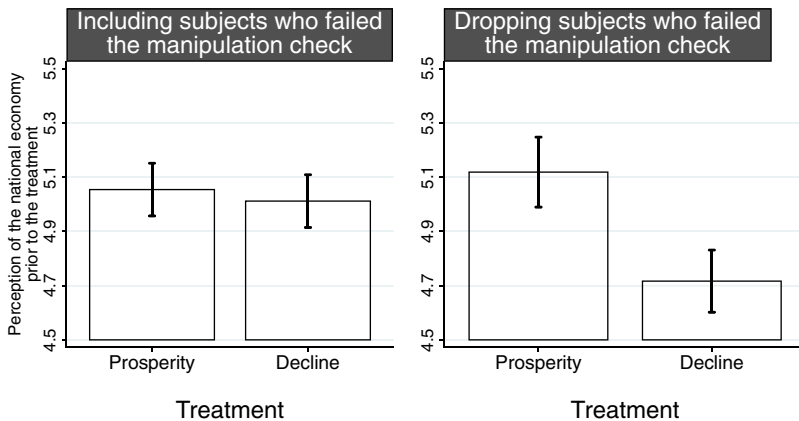


Figure 3

Average perception of the national economy prior to the experiment depending on the experimental treatment (with 95% Confidence Interval).

characteristics prior to the experiment. In contrast, when excluding subjects who failed the manipulation check, there is a significant difference across treatments in the perception of the national economy prior to the survey experiment, $t(2366) = 4.5688$, $p < .0001$. It is impossible that the experimental treatment had a causal effect on responses to a question asked earlier in the survey. Thus, this reflects a selection effect emerging from dropping subjects who failed the manipulation check.

Suppose that we want to test the effect of the treatment on a dependent variable. After the treatment, the subjects' level of nostalgia was assessed. Subjects indicated on a 5-point scale from "strongly disagree" to "strongly agree" to what extent they agreed that "the society used to be a much better place". As shown by regression models in Table 1, when dropping subjects who failed the manipulation check (model 1), one would conclude that the *decline* treatment had a significant positive effect on nostalgia. Yet, this effect is reduced when controlling for the preexisting perception of the economy (model 2). Since it is impossible to measure all potential preexisting characteristics of subjects selected through the manipulation check, the best option is to avoid dropping subjects (model 3). This decision is conservative: it is likely to greatly increase the noise in the data and reduce effect sizes – with increased risk of Type II error. In model 3, the effect of the treatment is no longer significant³. Nonetheless, this is the only way to ensure that, if some treatment effect is observed, it is of genuinely causal nature.

The present study does not advocate against the inclusion of posttreatment manipulation checks. These can be informative tools – especially in the development phase of experiments – to assess the degree of attention and comprehension of the treatment in the given type of sample. In the present study, the manipulation check reveals that a very large fraction of subjects were inattentive. One possibility is

³Note that results are essentially unchanged when including country fixed effects (see the online appendix).

Table 1
Results from linear regression models of the level of nostalgia

	(1)	(2)	(3)
	Dropping subjects who failed the manipulation check		All subjects
Decline treatment	0.164***	0.114**	0.0603
	(0.0433)	(0.0425)	(0.0337)
Perception of the economy prior to the experiment		-0.138***	
		(0.00989)	
Constant	3.344***	4.054***	3.426***
	(0.0323)	(0.0597)	(0.0237)
Observations	2395	2311	3797

Standard errors in parentheses.
* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

that the pretreatment question about the subjects’ perception of the national economy induced subjects to disregard the content of the experimental vignette. This would explain the high rate of failure in the responses to the manipulation check and their close correlation with subjects’ preexisting beliefs about the economy. To test for this, it would be necessary to have a control group for whom the initial question was not included. However, this limitation does not affect the overall conclusion that exclusion based on posttreatment manipulation checks must be avoided.

What are then the options available to researchers? A first option highlighted by the literature is to include pretreatment questions to gauge subjects’ attentiveness. A common tool is instructional manipulation checks – or “screeners”. Screeners are similar to classic survey questions but ask participants to ignore the standard response format and instead provide a confirmation that they have read the instruction Berinsky et al. (2014); Oppenheimer et al. (2009). One disadvantage is that screeners may induce subjects to think that researchers want to trap them, which alters their responses to subsequent questions (Hauser and Schwarz 2015). Alternatively, Kane et al. (2020) recently proposed ready-to-use “mock vignettes”. A mock vignette mimics common kind of descriptive content in political science experiments but appear before the researcher’s treatment. All subjects read the same vignette and must then answer factual questions about it, allowing to check for subjects’ attentiveness.

The latter tools come with the cost of sacrificing survey space. Another alternative is to rely on timers as a proxy to identify inattentive subjects. Various studies highlight that subjects with short response times are generally less attentive (see for instance Börger 2016; Wood et al. 2017)⁴. Read et al. (2021) designed a method to

⁴Our study includes a measure of the overall duration of the survey, which confirms this. The share of subjects who failed the manipulation check question is significantly higher among subjects who spent relatively less time in the survey (see the online appendix).

identify inattentive subjects based on multiple question timers. Their method does not induce posttreatment selection bias when computed on question timers before the treatment. Besides, it allows to identify slow but nonetheless inattentive subjects.

Depending on the space available in survey, researchers may use these methods to perform analyses on sub-sample(s) of attentive subjects, in order to limit the risk of Type II error without inducing posttreatment bias. These measures of attentiveness may correlate with politically relevant variable, such as age, race, and education (see Berinsky et al. 2014; Kane et al. 2020). Thus, restricting analyses to attentive subjects comes with the risk of drawing conclusions that are not representative of the population. To mitigate this risk, the best practice should be to report estimates of treatment effects based on both the overall sample and subsample(s) of attentive subjects.

Data Availability Statement. This research is part of the Popeuropa project supported by IDEX Université Grenoble Alpes (IRS 2017-2018), Sciences Po Grenoble, and Pacte laboratory. The data, code, and any additional materials required to replicate all analyses in this article are available at the Journal of Experimental Political Science Dataverse within the Harvard Dataverse Network at <https://doi.org/10.7910/DVN/7DXBGG>.

Acknowledgements. I would like to thank two anonymous reviewers for their insightful comments that greatly contributed to improving the manuscript. I am grateful to Antoine Machut and the team of the Vendredis Quanti network for early discussions about this study. I would also like to thank Francis Varaine for his comments on the manuscript.

Conflicts of Interest. I acknowledge that I have no conflicts of interest or potential conflicts of interest with regard to the submitted work.

Ethics Statement. No IRB approval was sought for this project. The study complies with the European General Data Protection Regulation relative to the protection of personal data and received the approval of the scientific committee of Sciences Po Grenoble. We obtained informed consent from all participants, who could choose not to answer any questions or withdraw from the study at any time. Compensation was delivered by the survey vendor. This research adheres to APSA's Principles and Guidance for Human Subjects Research. The experiment included short written abstracts of fictional news media prospects about the economy. Given the variety of economic prospects commonly available in public news media, it was considered that no significant deception was induced. Section 1 in the online appendix details the experimental procedure employed.

References

- Aronow, P. M., J. Baron, and L. Pinson 2019. A Note on Dropping Experimental Subjects Who Fail a Manipulation Check. *Political Analysis* 27, 572–89.
- Berinsky, A. J., M. F. Margolis, and M. W. Sances 2014. Separating the Shirkers from the Workers? Making Sure Respondents Pay Attention on Self-administered Surveys. *American Journal of Political Science* 58, 739–53.
- Börger, T. 2016. Are Fast Responses More Random? Testing the Effect of Response Time on Scale in an Online Choice Experiment. *Environmental and Resource Economics* 65, 389–413.
- Hauser, D. J. and N. Schwarz 2015. It's a Trap! Instructional Manipulation Checks Prompt Systematic Thinking on “Tricky” Tasks. *Sage Open* 5, 2158244015584617.
- Kane, J. V. and J. Barabas 2019. No Harm in Checking: Using Factual Manipulation Checks to Assess Attentiveness in Experiments. *American Journal of Political Science* 63, 234–49.
- Kane, J. V., Y. R. Velez, and J. Barabas 2020. Analyze the Attentive & Bypass Bias: Mock Vignette Checks in Survey Experiments. *APSA Preprints*. doi: [10.33774/apsa-2020-9672-v2](https://doi.org/10.33774/apsa-2020-9672-v2).

- Montgomery, J. M., B. Nyhan, and M. Torres** 2018. How Conditioning on Posttreatment Variables can Ruin Your Experiment and What to do About it. *American Journal of Political Science* 62, 760–75.
- Mutz, D. C. and R. Pemantle** 2015. Standards for Experimental Research: Encouraging a Better Understanding of Experimental Methods. *Journal of Experimental Political Science* 2, 192–215.
- Oppenheimer, D. M., T. Meyvis, and N. Davidenko** 2009. Instructional Manipulation Checks: Detecting Satisficing to Increase Statistical Power. *Journal of Experimental Social Psychology* 45, 867–872.
- Read, B., L. Wolters, and A. J. Berinsky** 2021. Racing the Clock: Using Response Time as a Proxy for Attentiveness on Self-Administered Surveys. *Political Analysis*, 1–20. doi: [10.1017/pan.2021.32](https://doi.org/10.1017/pan.2021.32).
- Sansone, C., C. C. Morf, and A. T. Panter** 2003. *The Sage Handbook of Methods in Social Psychology*. Thousand Oaks: Sage Publications.
- Stenner, K.** 2005. *The Authoritarian Dynamic*. Cambridge: Cambridge University Press.
- Varaine, S.** 2022. Replication Data for: How dropping subjects who failed manipulation checks can bias your experimental results. An illustrative case. *Harvard Dataverse*. doi: [10.7910/DVN/7DXBGG](https://doi.org/10.7910/DVN/7DXBGG).
- Wood, D., P. D. Harms, G. H. Lowman, and J. A. DeSimone** 2017. Response Speed and Response Consistency as Mutually Validating Indicators of Data Quality in Online Samples. *Social Psychological and Personality Science* 8, 454–464.

Cite this article: Varaine S (2023). How Dropping Subjects Who Failed Manipulation Checks Can Bias Your Results: An Illustrative Case. *Journal of Experimental Political Science* 10, 299–305. <https://doi.org/10.1017/XPS.2022.28>