

# Why overfitting is bad for science

## LESSONS FROM PSYCHOLOGY

13/02/2020 | eResearchNZ

---

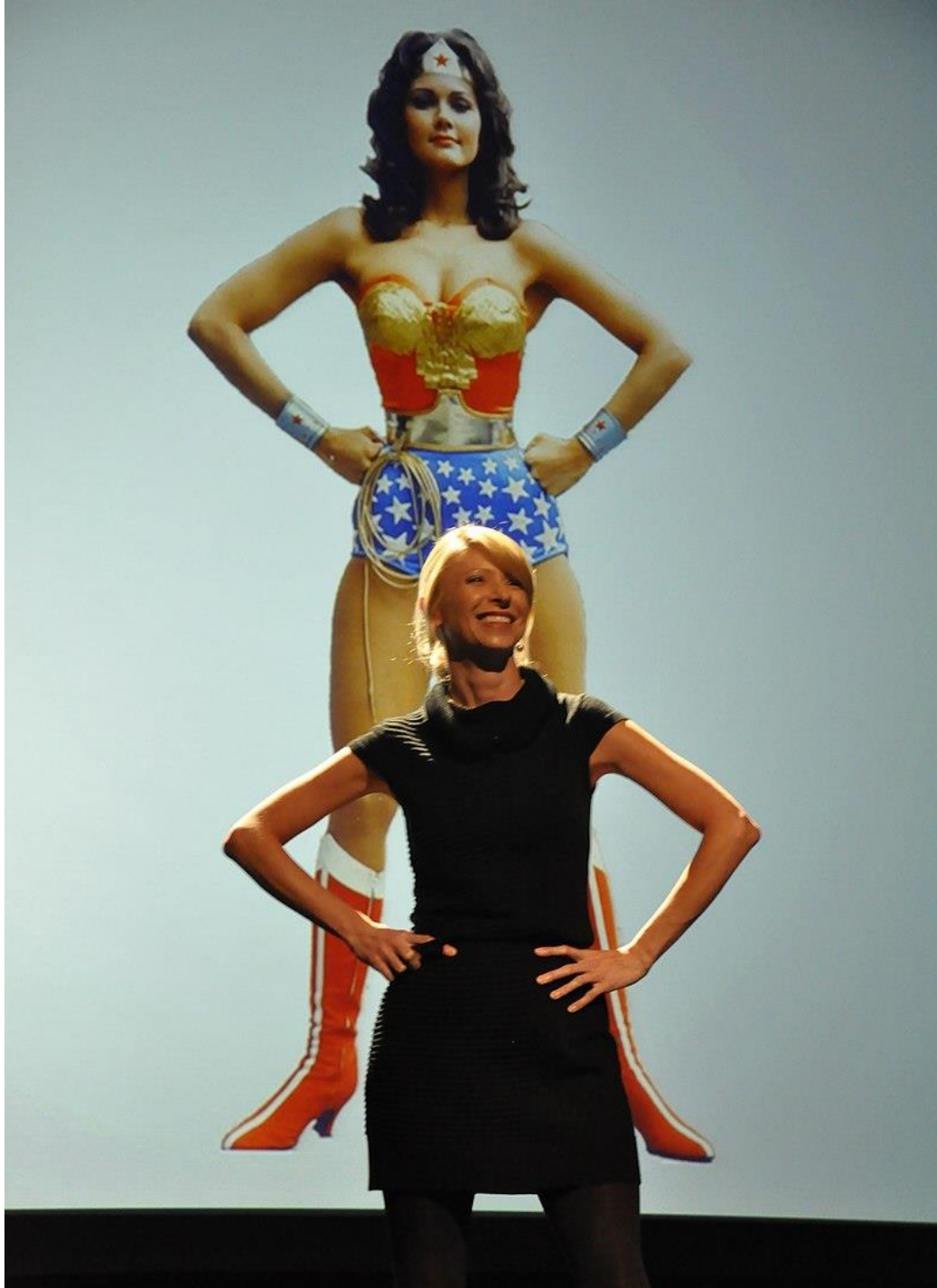
Adam Bartonicek | PhD student | Department of Psychology

Narun Pat | Lecturer | Department of Psychology

Tamlin Conner | Associate Professor | Department of Psychology



@BartonicekAdam



- "Power posing: Brief nonverbal displays affect neuroendocrine levels and risk tolerance."
- (Carney, Cuddy, and Yap, 2010)



Powerposing did not replicate!



# Many psychology papers fail replication test

An effort to repeat 100 studies yields sobering results, but many researchers are positive about the process

By John Bohannon

The largest effort yet to replicate psychology studies has yielded both good and bad news. On the down side, of the 100 prominent papers analyzed, only 39% could be replicated unambiguously, as a group of 270 researchers describes on page 943. On the up side, despite the sobering results, the effort seems to have drawn little of the animosity that greeted a similar replication effort last year (*Science*, 23 May 2014, p. 788). This time around, many of the original authors are praising the replications as a useful addition to their own research.

## RESEARCH ARTICLE

### PSYCHOLOGY

## Estimating the reproducibility of psychological science

Open Science Collaboration<sup>\*†</sup>

Replication effects were half the magnitude of original effects, representing a substantial decline. Ninety-seven percent of original studies had statistically significant results. Thirty-six percent of replications had statistically significant results; 47% of original effect sizes were in the 95% confidence interval of the replication effect size; 39% of effects were subjectively rated to have replicated the original result; and if no bias in original results is assumed, combining original and replication results left 68% with statistically significant effects.

Perspectives on Psychological Science  
2016, Vol. 11(4) 546–573  
© The Author(s) 2016  
Reprints and permissions:  
sagepub.com/journalsPermissions.nav  
DOI: 10.1177/1745691616652873  
pps.sagepub.com

### A Multilab Preregistered Replication of the Ego-Depletion Effect

M. S. Hagger,<sup>\*</sup> N. L. D. Chatzisarantis,<sup>\*</sup> H. Alberts, C. O. Anggono, C. Batailler, A. R. Birt, R. Brand, M. J. Brandt, G. Brewer, S. Bruyneel, D. P. Calvillo, W. K. Campbell, P. R. Cannon, M. Carlucci, N. P. Carruth, T. Cheung, A. Crowell, D. T. D. De Ridder, S. Dewitte, M. Elson, J. R. Evans, B. A. Fay, B. M. Fennis, A. Finley, Z. Francis, E. Heise, H. Hoemann, M. Inzlicht, S. L. Koole, L. Koppel, F. Kroese, F. Lange, K. Lau, B. P. Lynch, C. Martijn, H. Merckelbach, N. V. Mills, A. Michirev, A. Miyake, A. E. Mosser, M. Muise, D. Muller, M. Muzi, D. Nalis, R. Nurwanti, H. Otgaar, M. C. Philipp, P. Primoceri, K. Rentzsch, L. Ringos, C. Schlankert, B. J. Schmeichel, S. F. Schoch, M. Schrama, A. Schütz, A. Stamos, G. Tinghög, J. Ullrich, M. vanDellen, S. Wimbarti, W. Wolff, C. Yussainy, O. Zerhouni, and M. Zwienerberg

<sup>\*</sup>Proposing authors

for the current Registered Replication Report of the ego-depletion effect. Multiple laboratories ( $k = 23$ , total  $N = 2,141$ ) conducted replications of a standardized ego-depletion protocol based on a sequential-task paradigm by Sripada et al. Meta-analysis of the studies revealed that the size of the ego-depletion effect was small with 95% confidence intervals (CIs) that encompassed zero ( $d = 0.04$ , 95% CI  $[-0.07, 0.15]$ ). We discuss implications of the findings for the ego-depletion effect and the resource depletion model of self-control.

## Thirty-six percent of replications had statistically significant results;

in the same direction as the original study for 13 (62%) studies, and the effect size of the replications is on average about 50% of the original effect size.

## Evaluating the replicability of social science experiments in *Nature* and *Science* between 2010 and 2015

Colin F. Camerer<sup>1,16</sup>, Anna Dreber<sup>2,16</sup>, Felix Holzmeister<sup>3,16</sup>, Teck-Hua Ho<sup>4,16</sup>, Jürgen Huber<sup>3,16</sup>, Magnus Johannesson<sup>5,16</sup>, Michael Kirchler<sup>3,5,16</sup>, Gideon Nave<sup>6,16</sup>, Brian A. Nosek<sup>7,8,16\*</sup>, Thomas Pfeiffer<sup>9,16</sup>, Adam Altmeld<sup>10,2</sup>, Nick Buttrick<sup>7,8</sup>, Taizan Chan<sup>10</sup>, Yiling Chen<sup>11</sup>, Eskil Forsell<sup>12</sup>, Anup Gampa<sup>7,8</sup>, Emma Heikensten<sup>1</sup>, Lily Hummer<sup>1</sup>, Taisuke Imai<sup>13</sup>, Siri Isaksson<sup>2</sup>, Dylan Manfredi<sup>16</sup>, Julia Rose<sup>2</sup>, Eric-Jan Wagenmakers<sup>14</sup> and Hang Wu<sup>15</sup>

We replicate 21 systematically selected experimental studies in the social sciences published in *Nature* and *Science* between 2010 and 2015<sup>16–36</sup>. The replications follow analysis plans reviewed by the original authors and pre-registered prior to the replications. The replications are high powered, with sample sizes on average about five times higher than in the original studies. We find a significant effect in the same direction as the original study for 13 (62%) studies, and the effect size of the replications is on average about 50% of the original effect size.

We find a significant effect in the same direction as the original study for 13 (62%) studies, and the effect size of the replications is on average about 50% of the original effect size.



### RESEARCH ARTICLE

## Failed Replication of Oxytocin Effects on Trust: The Envelope Task Case

Anthony Lane<sup>1\*</sup>, Moira Mikolajczak<sup>1</sup>, Evelynne Treinen<sup>2</sup>, Dana Samson<sup>1</sup>, Olivier Corneille<sup>1</sup>, Philippe de Timary<sup>3</sup>, Olivier Luminet<sup>1</sup>

In this paper we present two failed replications of this effect, despite sufficient power to replicate the original large effect. The non-significant results of these two failed replications clearly exclude a large effect of OT on trust in this paradigm but are compatible with either a null effect of OT on trust, or a small effect, undetectable with small sample size ( $N = 95$  and  $61$  in Study 1 and 2, respectively).

Only psychology in  
hot water?



"It's just a simple Rorschach ink-blot test, Mr. Bromwell, so just calm down and tell me what each one suggests to you."

# Why Most Published Research Findings Are False

**John P. A. Ioannidis**



PLOS Medicine | [www.plosmedicine.org](http://www.plosmedicine.org)

August 2005 | Volume 2 | Issue 8 | e124



# Power failure: why small sample size undermines the reliability of neuroscience

Katherine S. Button<sup>1,2</sup>, John P. A. Ioannidis<sup>3</sup>, Claire Mokrysz<sup>1</sup>, Brian A. Nosek<sup>4</sup>, Jonathan Flint<sup>5</sup>, Emma S. J. Robinson<sup>6</sup> and Marcus R. Munafò<sup>1</sup>

NATURE REVIEWS | NEUROSCIENCE

VOLUME 14 | MAY 2013 | 1

A quarter of most-cited clinical trials and 5/6 most-cited epidemiological studies were either fully contradicted or found to have exaggerated results<sup>2</sup>

## Why Most Discovered True Associations Are Inflated

John P. A. Ioannidis

Epidemiology • Volume 19, Number 5, September 2008

TABLE 1. Selected Evaluations Suggesting That Early Discovered Effects Are Inflated

Research Field	Theoretical Work or Empirical Evidence and References
Highly cited clinical research	A quarter of most-cited clinical trials and 5/6 most-cited epidemiological studies were either fully contradicted or found to have exaggerated results <sup>2</sup>
Early stopped clinical trials	Early stopping results in inflated effects in theory <sup>3,4</sup> and shown also in practice <sup>5</sup>
Clinical trials of mental health interventions	More likely for effect sizes of pharmacotherapies to diminish than to increase over time <sup>6</sup>
Clinical trials on heart failure interventions	"Regression to the truth" in phase III trials for interventions with early promising results <sup>7</sup>
Clinical trials on diverse interventions	Effectiveness shown to fade over time <sup>8</sup>
Multiple meta-analyses on effectiveness	Eleven independent meta-analyses on acetylcysteine show decreasing effects over time <sup>9</sup>
Epidemiologic associations	Expected to be inflated in multiple testing with significance threshold; empirical demonstration for occupational carcinogens <sup>10</sup>
Pharmacoepidemiology	"Phantom ship" associations that do not stand upon further evaluation <sup>11</sup>
Gene-disease associations	Several empirical evaluations showing dissipation of effect sizes over time <sup>12-15</sup>
Linkage studies in humans	Theory anticipates large upward bias ("winner's curse") in effects of discovered loci <sup>16-18</sup>
Genetic traits in experimental crosses	As above (actually literature on the "Beavis effect" precedes literature on humans) <sup>19-22</sup>
Genome-wide associations	Large winner's curse anticipated for discovered effects in underpowered conditions <sup>23,24</sup>
Ecology and evolution	Empirical demonstration that relationships fade over time <sup>25,26</sup>
Psychology	Replication studies in psychology failing to confirm true effects because the new studies were underpowered due to reliance on the estimate of effect from the original positive study <sup>27</sup>
Early repeated data peaking in general	Simulations to model inflation of effects with repeated data peaking <sup>28</sup>
Prognostic models	Overestimated prognostic performance with stepwise selection of variables based on significance thresholds <sup>29-32</sup>
Regression models in general	Exaggerated effects (coefficients) with stepwise selection based on significance thresholds and small datasets <sup>32-34</sup> ; may correct substantially if a very lenient alpha = 0.20 is used for selection <sup>34</sup> [thus having enough power]

## Believe it or not: how much can we rely on published data on potential drug targets?

Florian Prinz, Thomas Schlange and Khusru Asadullah

NATURE REVIEWS | DRUG DISCOVERY

We received input from 23 scientists (heads of laboratories) and collected data from 67 projects, most of them (47) from the field of oncology. This analysis revealed that only in ~20–25% of the projects were the relevant published data completely in line with our in-house findings (FIG. 1c). In almost two-thirds of the projects, there were inconsistencies between published data and in-house data that either considerably prolonged the duration of the target validation process or, in most cases, resulted in termination of the projects because the evidence that was generated for the therapeutic hypothesis was insufficient to justify further investments into these projects.

JAMA Published online August 23, 2018

John P. A. Ioannidis, MD, DSc  
Stanford Prevention Research Center and Meta-Research Innovation Center at Stanford (METRICS), Stanford University, Stanford, California.

VIEWPOINT

## The Challenge of Reforming Nutritional Epidemiologic Research

In recent updated meta-analyses of prospective cohort studies, almost all foods revealed statistically significant associations with mortality risk.<sup>1</sup>

## Raise standards for preclinical cancer research

C. Glenn Begley and Lee M. Ellis propose how methods, publications and incentives must change if patients are to benefit.

29 MARCH 2012 | VOL 483 | NATURE | 531

Although we in the cancer field hoped that this would lead to more effective drugs, historically, our ability to translate cancer research to clinical success has been remarkably low<sup>1</sup>.

REPRODUCIBILITY OF RESEARCH FINDINGS			
Preclinical research generates many secondary publications, even when results cannot be reproduced.			
Journal impact factor	Number of articles	Mean number of citations of non-reproduced articles*	Mean number of citations of reproduced articles
>20	21	248 (range 3–800)	231 (range 82–519)
5–19	32	169 (range 6–1,909)	13 (range 3–24)

\*Results from ten-year retrospective analysis of experiments performed prospectively. The term 'non-reproduced' was assigned on the basis of findings not being sufficiently robust to drive a drug development programme.  
\*Source of citations: Google Scholar, May 2011.



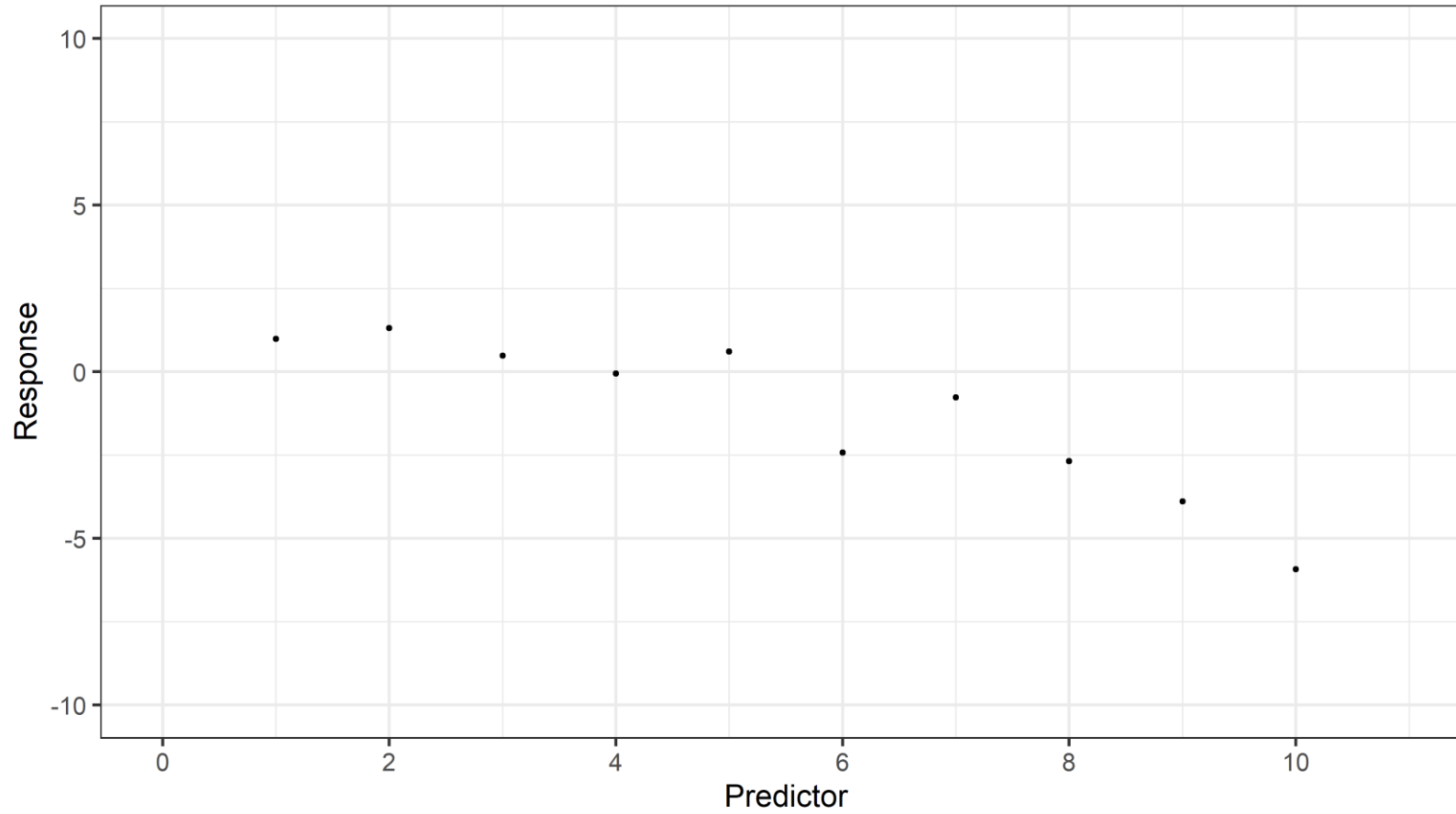
# Why are most research findings false?

- Low power (small N)
- Publication bias

# Why are most research findings false?

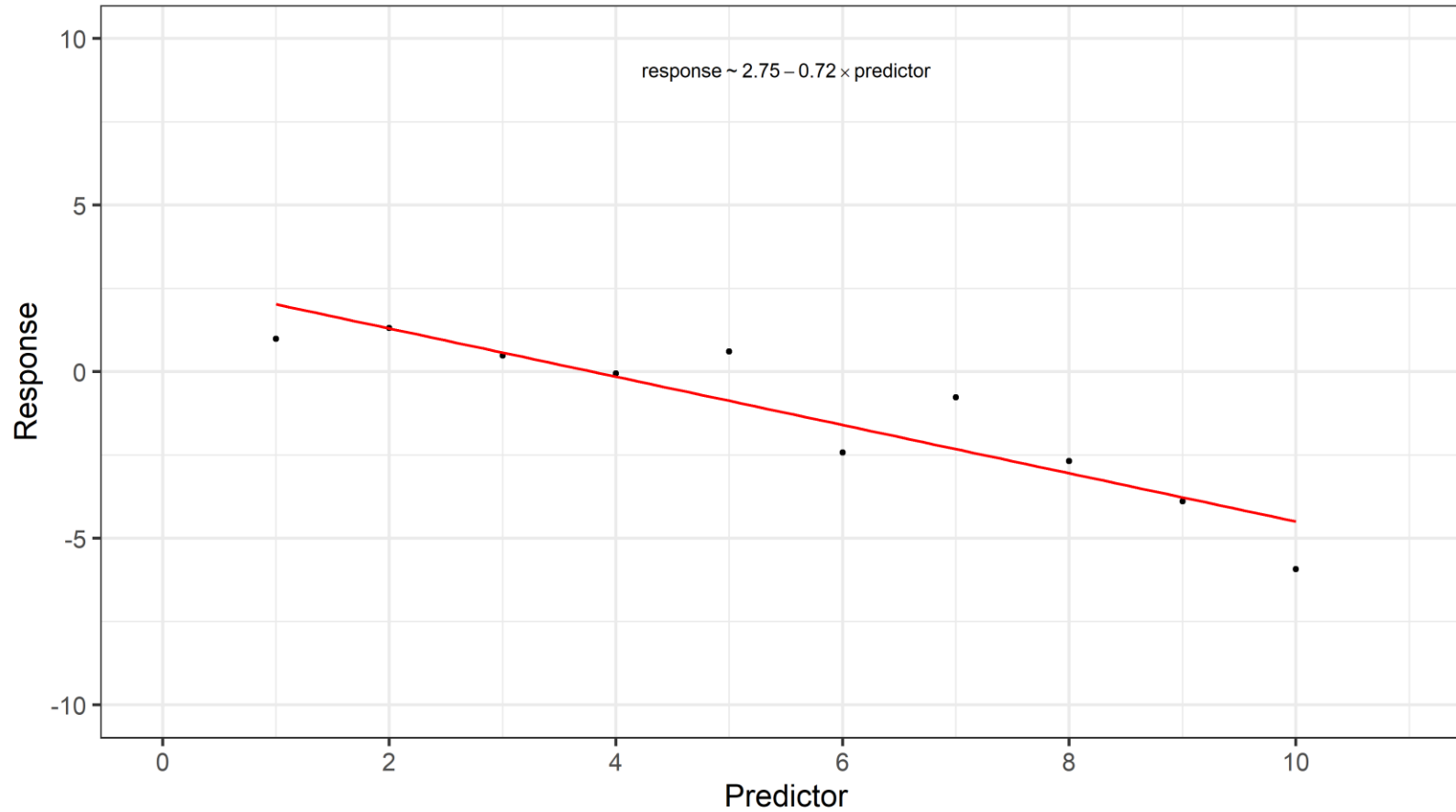
- Questionable research practices:
  - P-hacking
  - HARKing (hypothesizing after results are known)
  - Early stopping
  - (see Wicherts et al., 2016)
- (Overfitting?)

# What is overfitting?

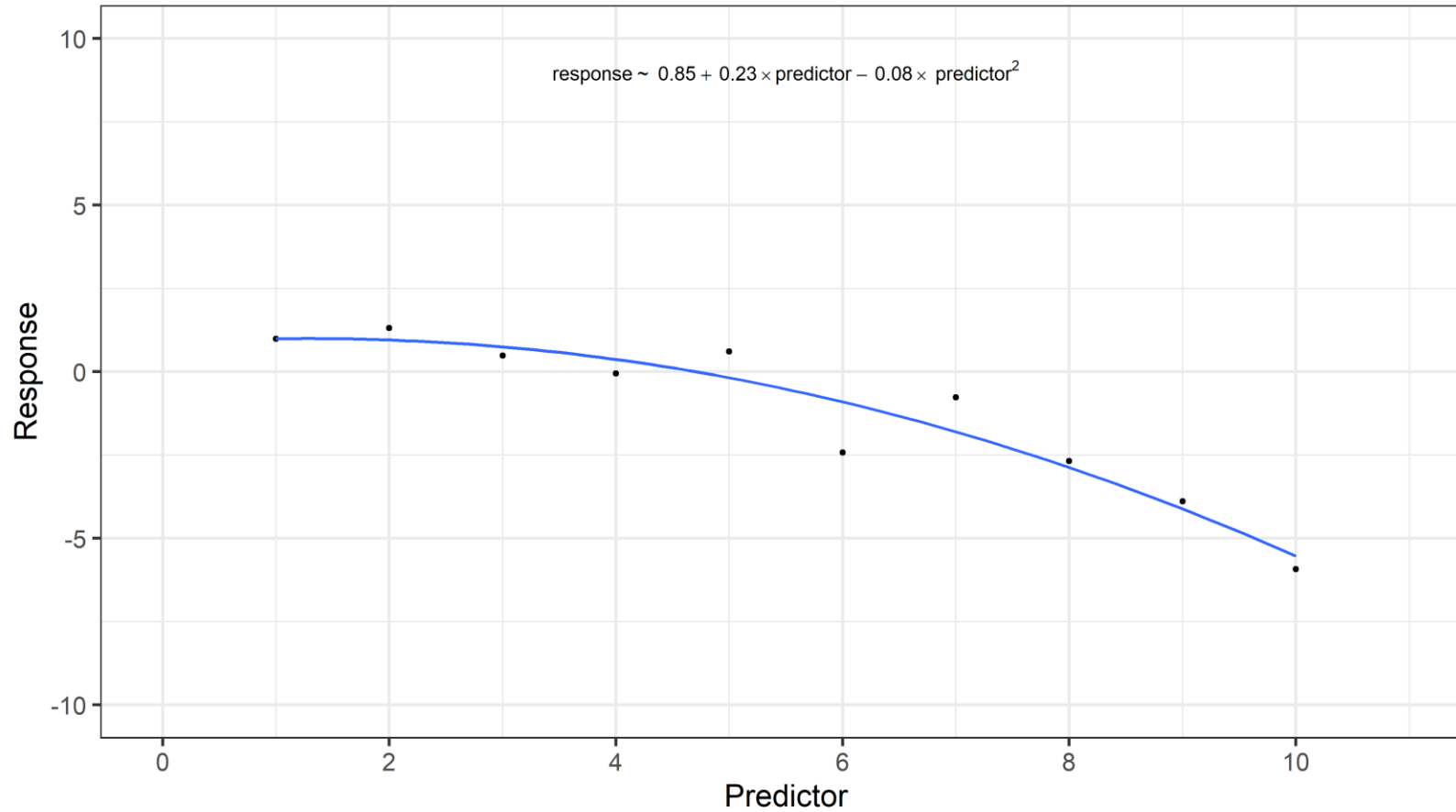




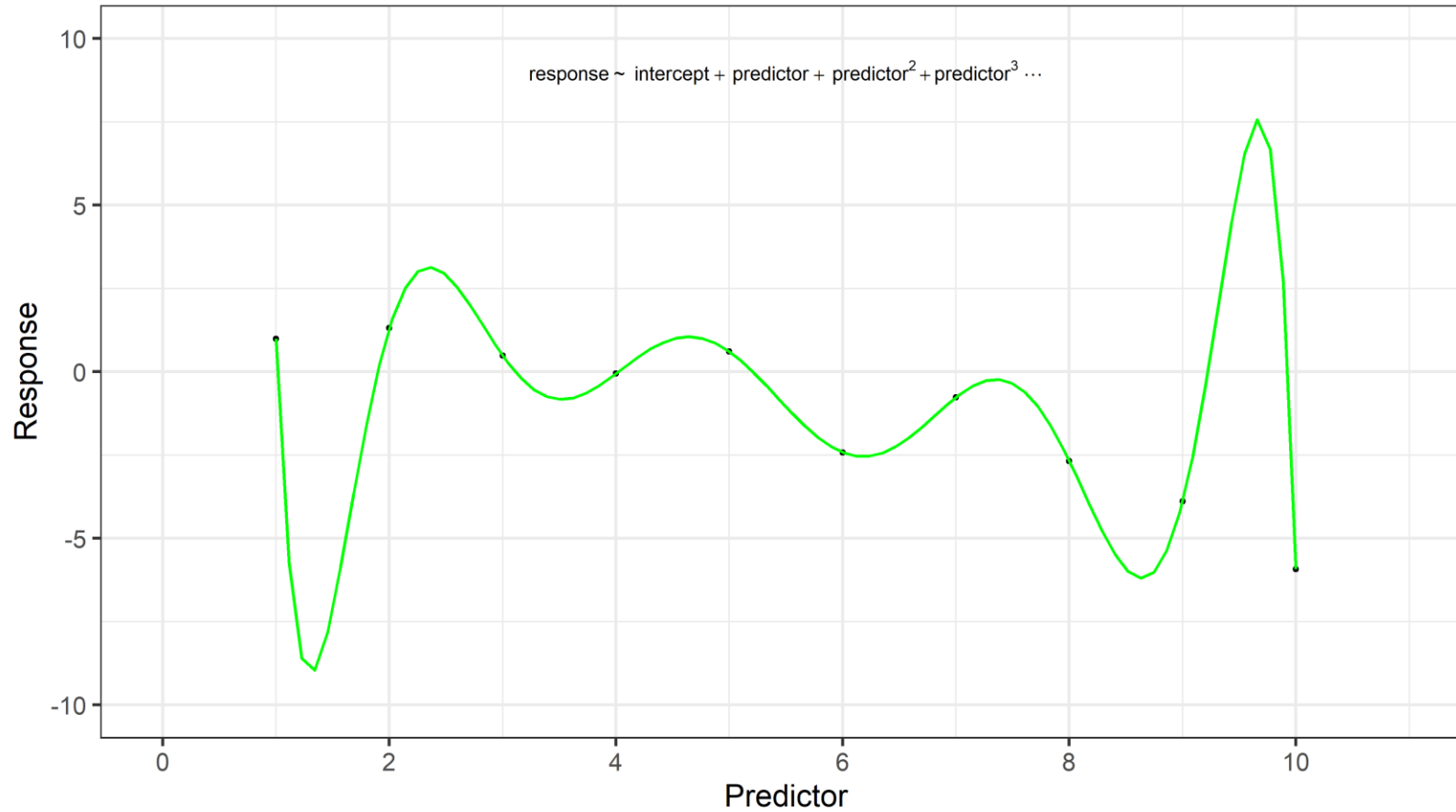
# What is overfitting?



# What is overfitting?



# What is overfitting?







# Overfitting

- More predictors  $\rightarrow$  flexible model  $\rightarrow$  better fit (e.g.  $R^2$ )
- Models can learn “too much” and fail to generalize to new samples
- Statistical significance  $\neq$  no overfitting

# Yarkoni & Westfall (2017)

- Replication crisis is the result of overfitting
- Ordinary least squares and maximum likelihood methods are vulnerable to overfitting
- Even worse with stepwise regression



# How to guard against overfitting?

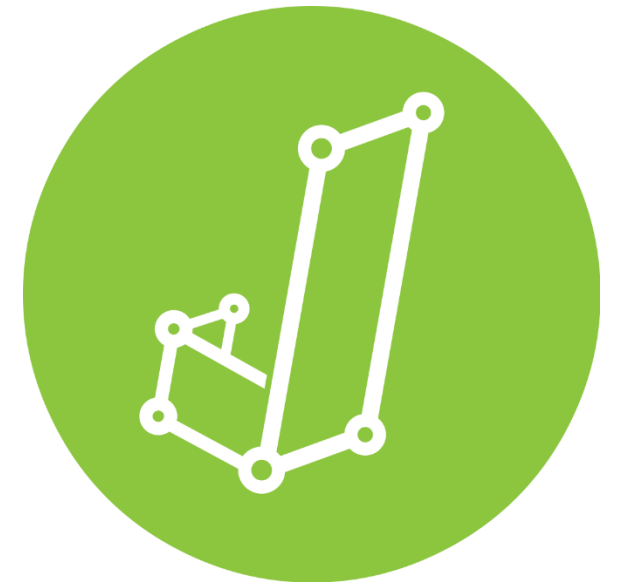


- Predictive power in addition to significance
- Machine learning methods:
  - Cross-validation
  - Regularization (ridge regression, LASSO, Bayesian priors with high probability density at 0, e.g. horseshoe)



# Open source ML software

- PredPsych (Koul, Becchio, & Cavallo, 2018)
  - User friendly R package with staple ML methods
- JASP machine learning module (2019)
  - JASP is a free, GUI statistical software
  - ML module as of 2019



# To summarize...

- Models can learn too much from the data and overfit
- Overfitting may lead to unreliable science and failed replications
- To avoid overfitting, we can use ML methods such as cross-validation and regularization

# Thank you!

Carney, D. R, A. J. Cuddy, and A. J. Yap (2010). "Power Posing: Brief Nonverbal Displays Affect Neuroendocrine Levels and Risk Tolerance". In: *Psychological Science* 21.10, pp. 1363-1368. ISSN: 14679280.

Wicherts, J. M, C. L. Veldkamp, H. E. Augusteijn, et al. (2016). "Degrees of freedom in planning, running, analyzing, and reporting psychological studies: A checklist to avoid P-hacking". In: *Frontiers in Psychology* 7.NOV, pp. 1-12. ISSN: 16641078.

Yarkoni, T. and J. Westfall (2017). "Choosing Prediction Over Explanation in Psychology: Lessons From Machine Learning". In: *Perspectives on Psychological Science* 12.6, pp. 1100-1122. ISSN: 17456924.