

# Criticism as asynchronous collaboration: An example from social science research\*

Andrew Gelman

11 Feb 2022

## 1. Introduction

Collaboration is essential to statistics: applied problems motivate the development and evaluation of new methods, which in turn allow new applied problems to be solved. This cycling between concerns of statistics and subject matter will only work in the presence of strong connections between statistics and applied fields.

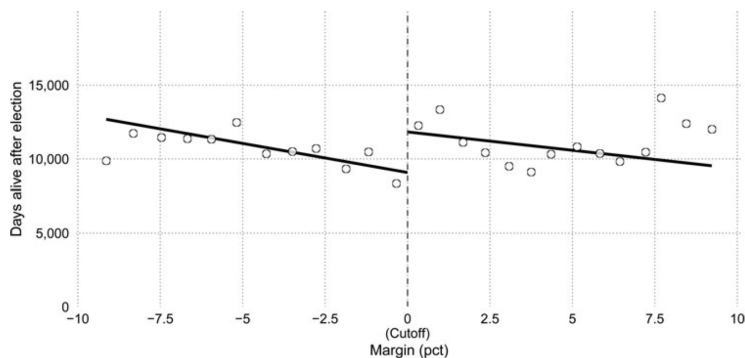
The present article illustrates a different sort of collaboration, performed implicitly between the authors of a published paper and later researchers who want to understand or use the published work. This is not a collaboration in the usual sense, as it does not require that the different research groups ever contact each other—but it is arguably the most important, and presumably the most common, form of research connection. Just about anyone who publishes a paper with a statistical result has a hope that various complete strangers will read it and take its findings seriously enough to attempt to understand and replicate them.

The most usual forms of scholarly engagement with published work are citation, appreciation, and refutation. But these do not fully capture the experience of reading and evaluating a paper, in particular the uncertainties involved in balancing theoretical and evidential claims.

Here we discuss a published paper in political science that made a claim—“politicians winning a close election live 5–10 years longer than candidates who lose”—that aroused my skepticism. My purpose here is not to formally rebut that paper (although we do present several arguments against it) but rather to use this as an example of how we, as consumers as well as producers of science, can engage with published work. Fortunately in this case the data from the published study were available to all for reanalysis.

## 2. Reactions to a published paper

Barfort, Klemmensen, and Larsen (2021) write: “we exploit a regression discontinuity design with unique data on the longevity of candidates for US gubernatorial office. The results show that politicians winning a close election live 5–10 years longer than candidates who lose.” Their analytic approach and results are well summarized by the following graph, which they labeled, “The causal effect of winning on longevity.”



\*We thank Erik Gahner Larsen and Bill Harris for helpful comments and the U.S. Office of Naval Research, Institute for Education Sciences, and Sloan Foundation for partial support of this work.

This regression discontinuity analysis makes use of the insight that, if treatment assignment (in this case, a candidate losing or winning the election) is entirely dependent on the running variable (in this case, the difference in vote share between the candidate), then there is no possibility of selection bias: in econometric or statistical jargon, the assignment is exogenous or ignorable conditional on this variable. The challenge of is that, by design, there is imbalance and zero overlap of treatment and control groups with respect to the running variable, hence the problem is typically framed in terms of statistical modeling and robustness, as the inference for the causal effect will depend crucially on the method used to adjust for this imbalance.

Despite the  $p < 0.05$  statistical significance of the discontinuity in the above graph, and despite the assurance of the authors that this finding is robust, no, I do not believe that winning a close election causes U.S. governors to live 5–10 years longer.

How can I say this? I will answer in five ways:

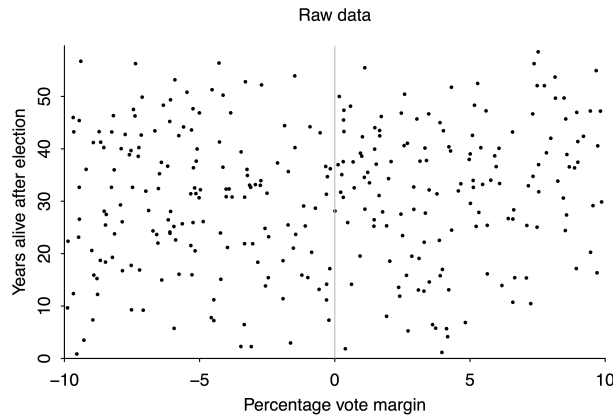
1. Common sense
2. Our experience as consumers of research
3. Statistical analysis
4. Statistical design
5. Sociology of science.

Our discussion of these five reasons will shed some light, we hope, on the ways in which authors and researchers interact.

**Common sense.** Five to ten years is a huge amount of lifetime. Even if you imagine dramatic causal sequences (for example, you lose an election and become depressed and slide into alcoholism and then drive your car into a tree), it’s hard to imagine such a huge effect. For statistical reasons discussed below, it is not surprising that this sort of study can come up with an effect size estimate that is implausibly large, even in the absence of any true effect of that magnitude. For now, though, I will just say that such a large effect, while not physically impossible, gives us reason for doubt.

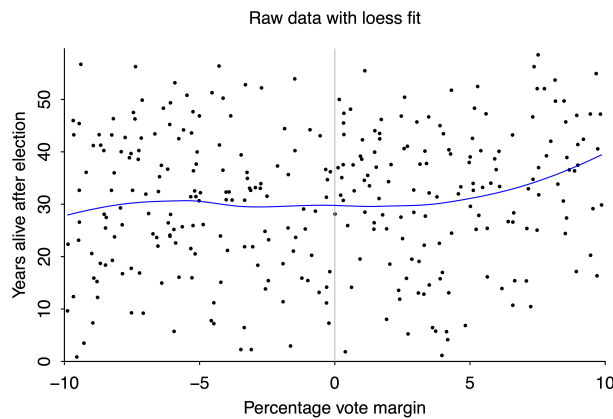
**Our experience as consumers of research.** The recent history of social and behavioral science is littered with published papers that made claims of implausibly large effects, supported by statistically significant comparisons and seemingly solid research methods, which did not hold up under replication or methodological scrutiny. Some examples that I happened to have looked into in some detail include the claim that beautiful parents were eight percentage points more likely to have girls, the claim that students at Cornell had extra-sensory perception, the claim that women were three times more likely to wear red or pink clothing during certain times of the month, the claim that single women were twenty percentage points more likely to support Barack Obama during certain times of the month, and the claim that political moderates perceived the shades of gray more accurately than extremists on the left and right. We’ve been burned before. We’ve been burned enough times that we realize we don’t have to follow the now-retired dictum of Kahneman (2011) that “you have no choice but to accept that the major conclusions of these studies are true.” For a particularly vivid example, we recommend the story of Nosek, Spies, and Motyl (2013) of a theoretically founded, statistically significant result they found in one of their research studies—which then did not show up in an attempted replication. It can be as easy to fool oneself as to fool others, and we know by now that confidence in the part of authors and publication in a respectable journal is not enough. It is not that nothing should be trusted; rather, the point here is that published results in the human sciences are often mistaken. In particular, common methods of statistical analyses routinely lead to overstatement of evidence.

**Statistical analysis.** If there truly is no such large effect that losing an election causing you to lose 5 to 10 years of life, then how could these researchers have found a comparison that was (a) statistically significant, and (b) robust to model perturbations? My quick answer is researcher degrees of freedom or forking paths (Simmons, Nelson, and Simonsohn, 2011, Gelman and Loken, 2014): the many different ways the analyses could've been done, contingent on data. The above graph might look compelling, but here's what the raw data look like:

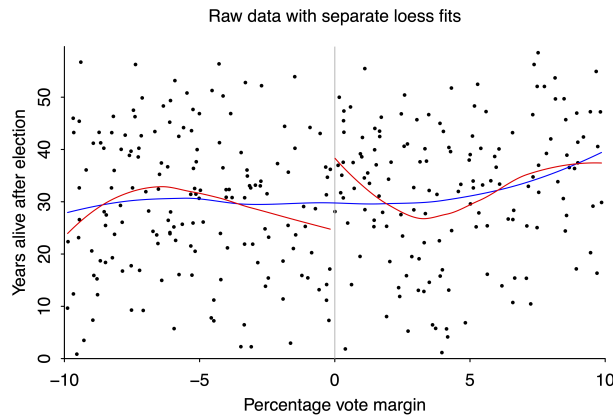


For ease of interpretation I've plotted the data in years rather than days.

Now let's throw a local linear smoother (loess) on there and see what we get:



And here are two loess fits, one for negative  $x$  and one for positive  $x$ :



We've seen this movie before (for example, Gelman and Zelizer, 2015). Fit a negative slope

right near the discontinuity, and then a positive jump is needed to make everything work out. The point is not that loess is the right thing to do here; the point is that this is what’s in these data.

The fit is noisy, and finding the discontinuity all depends on there being this strong negative relation between future lifespan and vote margin in this one election—but just for vote margins in this  $\pm 5$  percentage point range. Without that negative slope, the discontinuity goes away.

At this point you might say, no, the authors actually fit a local linear regression, so we can’t blame the curve, and that their results were robust . . . OK, we’ll get to that. My first point here is that the data are super-noisy, and fitting different models to these data will give you much different results. Again, remember that it makes sense that the data are noisy—there’s *no good reason at all* to expect a strong relationship between vote margin in an election and the number of years that someone will live afterward. Indeed, from any usual way of looking at things, it’s ludicrous to think that a candidate’s life expectancy is:

30 years if he loses an election by 5 percentage points

25 years if he loses narrowly

35 years if he wins narrowly

30 years if he wins by 5 percentage points.

It’s a lot more believable that this variation is just noise, some artifact of the few hundred cases in this dataset, than that it represents some general truth about elections, or even about elections for governor.

And here’s a linear regression. We’ll have more regressions in the next section; here’s one including all the elections between 1945 and 2012, excluding all missing data, and counting the first race for each candidate who ran for governor multiple times:

```
                coef.est coef.se
(Intercept)    78.60    4.05
won             2.39    2.44
age            -0.98    0.08
decades_since_1950 -0.21  0.51
margin         -0.11    0.22
---
n = 311, k = 5
residual sd = 10.73, R-Squared = 0.35
```

The estimated effect is 2.4 years with a standard error of 2.4 years, i.e., consistent with noise.

But what about the robustness as reported in the published article? My answer to that is, first, the result is not so robust, as is indicated by the above graph and regression and demonstrated further in the next section—and I wasn’t trying to make the result go away, I was just trying to replicate what they were doing in that paper,—and, second, as Simonsohn (2016) explains, “Robustness checks involve reporting alternative specifications that test the same hypothesis. Because the problem is with the hypothesis, the problem is not addressed with robustness checks.” Simonsohn illustrates that point with an amusing example in which he roots through the Generalized Social Survey to find a (spurious) relationship between the response to a horoscope-reading question and the serial number that was randomly assigned to survey respondents. His statistically significant result survives a reasonable-looking robustness check.

To return to the general theme of remote collaboration: the authors of a published paper supply robustness checks as part of an asynchronous dialogue with readers who might be skeptical of the claims. Alternative analyses are part of any good statistical study, but generally the point should be to explore and understand the limitations of one’s conclusions, not to rule out alternative explanations.

**Statistical design.** Another way to understand my skepticism is to consider the design of this sort of study study. A priori we might consider an effect size of one additional year of life to be large, and on the border of plausibility. But this study has essentially no power to detect effects this small! You can see that from the standard errors on the regression. If an estimate of 5–10 years is two or three standard errors from zero, than an effect of 1 year, or even 2 years, is statistically undetectable. So the study is really set up only to catch artifacts or noise. This is what Button et al. (2012) call “power failure.”

**Sociology of science.** If you are a social scientist, statistical methods should be your servant, not your master. It’s tempting to say that the authors of the paper in question followed the rules of regression discontinuity analysis and that it’s not fair to pick on them. Indeed, in other work, we have criticized the use of high-degree polynomials in discontinuity analyses (Gelman and Imbens, 2019), so you might even say that by using linear regression, the authors were heeding my own advice! But the point of social science is not to follow rules, it’s to gain understanding and make decisions. When following the rules gives silly results, it’s time to look more carefully at the rules and think hard about what went wrong. That’s how many in the field of psychology responded (appropriately, in my opinion) to the replication crisis of the 2010s.

Yes, it’s possible that I’m wrong. Perhaps losing a close election really did kill these candidates for governor. It’s possible. But I don’t think so. I think this paper is just another example of statistical rule-following that’s out of control.

### 3. Working with the data

I followed the link at the published article and downloaded the data. The code didn’t quite run as is—the R code required a .csv file but all I could find was a .tab file, so I changed the code accordingly. Then when I ran the next bit of the code:

```
if (sha1(df_rdd) != "300cc29bbeed2b630016c9bd2c8ef958dcc1b45d"){  
  error("Wrong data file loaded or data has been changed!") }
```

I got an error:

```
Error in error("Wrong data file loaded or data has been changed!") :  
  could not find function "error"
```

So I checked by typing:

```
sha1(df_rdd)
```

and got this:

```
"e26cc6acbed7a7144cd2886eb997d4ae262cf400"
```

So maybe the data did get changed! I have no idea. But I was able to reproduce the graphs so things were probably OK. I include these details not because they are important in themselves but as a small example of the operational details that rarely appear in print.

I also noticed some data problems, such as cases where the politician’s death date came decades before his election. There were a bunch of these, for example:

22	massachusetts	1867	adams	john	quincy	d	1826-07-04	no	challenger	-16.568840193044238
23	massachusetts	1868	adams	john	quincy	d	1826-07-04	no	challenger	-35.240317933127585
24	massachusetts	1869	adams	john	quincy	d	1826-07-04	no	challenger	-18.720612619251682
25	massachusetts	1870	adams	john	quincy	d	1826-07-04	no	challenger	-24.07333754454921
26	massachusetts	1871	adams	john	quincy	d	1826-07-04	no	challenger	-22.306152017842322
27	new hampshire	1848	adams	sherman	na	r	1886-10-27	no	challenger	4.97021484154476

I'm not trying to say that this is a horrible dataset. This is just what happens when you try to replicate someone's analyses. Problems come up.

Once I started to go through the data, I realized there are a lot of analysis choices. The article in question focuses on details of the discontinuity analysis, but that's really the least of it, considering that we have no real reason to think that the vote margin should be predictive of longevity. The most obvious thing is that the best predictor of how long you'll live in the future is ... how long you've lived so far. So I included current age as a linear predictor in the model. It then would be natural to interact the win/loss indicator with age. That to me makes a lot more sense than interacting with the vote margin. In general, it makes sense to interact the treatment with the most important predictors in your data.

Another issue is what to do with candidates who are still living. The published analysis just discarded them. It would be more natural, I think, to include such data using survival analysis. I didn't do this, though, as it would take additional work. I guess that was the motivation of the authors, too. No shame in that—I cut corners with missing data all the time—but it is yet another researcher degree of freedom.

Another big concern is candidates who run in multiple elections. It's not appropriate to use a regression model assuming independent errors when you have the same outcome for many cases. To keep things simple, I just kept the first election in the dataset for each candidate, but that's not the only choice; instead you might, for example, use the average vote margin for all the elections where the candidate ran.

Finally, for the analysis itself: it seems that the published regressions were performed using an automatic regression discontinuity package, `rdrobust`. I'm sure this is fine for what it is, but, again, I think that in this observational study the vote margin is not the most important thing to adjust for. Sure, you want to adjust for it, as there's no overlap on this variable, but we simply don't have enough cases right near the margin (the zone where we might legitimately approximate win or loss as a randomly applied treatment) to rely on the discontinuity alone. 400 or so cases is just not enough to estimate a realistic effect size here. To keep some control over the analyses, I just fit some simple regressions. It would be possible to include local bandwidths etc. if that were desired, but, again, that's not really the point. The discontinuity assignment is relevant to our analysis, but it's not the only thing going on, and we have to keep our eye on the big picture if we want to seriously estimate the treatment effect.

I then did some manipulations to check the dates in the file and clean the data, counting each candidate just once. My code was ugly. I did check the results a bit, but I wouldn't be surprised if I introduced a few bugs. I created a `decades_since_1950` variable as I had the idea that longevity might be increasing, and I put it in decades rather than years to get a more interpretable coefficient. I restricted the data to 1945–2012 and to candidates who were no longer alive at this time because that's what was done in the paper, and I considered election margins of less than 10 percentage points because that's what they showed in their graph, and also this did seem like a reasonable boundary for close elections that could've gone either way (so that we could consider it as a randomly assigned treatment).

Then I fit some regressions, with each time the outcome being the number of additional years lived by the candidate after the election:

```
coef.est coef.se
```

```

(Intercept)      78.60    4.05
won              2.39    2.44
age             -0.98    0.08
decades_since_1950 -0.21    0.51
margin          -0.11    0.22
---
n = 311, k = 5
residual sd = 10.73, R-Squared = 0.35

```

```

              coef.est coef.se
(Intercept)  78.81    5.41
won          1.92    8.25
age         -0.98    0.10
decades_since_1950 -0.21    0.51
margin      -0.11    0.22
won:age      0.01    0.16
---
n = 311, k = 6
residual sd = 10.75, R-Squared = 0.35

```

I was going to worry about how to interpret the coefficient of the treatment (“won” = winning the election) in the second model, because of the interaction with age. But the coefficient for that interaction is so small and so noisy that let’s just forget about it and go up to the earlier regression. The estimated treatment effect is 2.4 years (not 5—10 years) and its standard error is also 2.4.

What about excluding `decades_since_1950`?

```

              coef.est coef.se
(Intercept)  78.62    4.05
won          2.36    2.44
age         -0.99    0.08
margin      -0.11    0.22
---
n = 311, k = 4
residual sd = 10.72, R-Squared = 0.35

```

That didn’t do much. We could exclude age, which I wouldn’t recommend given how strong a predictor it is, but let’s see:

```

              coef.est coef.se
(Intercept)  30.14    1.75
won          1.70    3.01
margin      0.10    0.27
---
n = 311, k = 3
residual sd = 13.25, R-Squared = 0.01

```

Now the estimate’s even smaller and noisier! We should’ve kept age in the model in any case. We could up the power by including all elections where the vote margin was within 20 points, which gives us:

```

              coef.est coef.se
(Intercept)  76.24    3.29
won          1.00    1.83
age         -0.93    0.06

```

```

decades_since_1950 -0.18    0.39
margin              0.02    0.09
---
n = 497, k = 5
residual sd = 10.83, R-Squared = 0.33

```

Now we have almost 500 cases, but we're still not seeing that large and statistically significant effect.

I'm not saying that my regressions are absolutely better than the ones in the published paper. I'm pretty sure they're better in some ways and worse in others. I just found it surprisingly difficult to reproduce their results using conventional approaches. So, sure, I believe they found what they found ... but I call the result fragile, not robust.

Still not sure what to try, I re-ran the analysis going back to keeping only the elections with margins under 10 points, but including the duplicates as separate observations in the regression:

```

                coef.est coef.se
(Intercept)    74.65    3.18
won             3.15    1.89
age            -0.91    0.06
decades_since_1950 -0.03  0.41
margin         -0.19    0.17
---
n = 499, k = 5
residual sd = 10.93, R-Squared = 0.33

```

Almost statistically significant (whatever that's supposed to mean) with a  $z$ -score of 1.7, but still not in that 5–10 year range. Also, as we've already discussed, it doesn't make sense to count a politician multiple times, as each person only has one life.

I thought I could get cute and remove `age` and `decades_since_1950` and maybe something like the published paper's result would appear, but no luck:

```

                coef.est coef.se
(Intercept)  28.45    1.32
won           2.84    2.30
margin       -0.12    0.20
---
n = 499, k = 3
residual sd = 13.32, R-Squared = 0.00

```

I tried a few other things but I couldn't figure out how to get that huge and statistically significant coefficient in the 5–10 year range. One clue came from reproducing their results using `rdrobust`:

Method	Coef.	Std. Err.	$z$	$P> z $	[ 95% C.I. ]
Conventional	2749.283	873.601	3.147	0.002	[1037.057 , 4461.509]
Robust	-	-	3.188	0.001	[1197.646 , 5020.823]

I tried re-running the basic linear model but just in the zone where the candidates were within 5 percentage points of winning or losing, which yielded:

```

                coef.est coef.se
(Intercept)    71.03    6.87
won             7.54    3.76

```



```

age                -0.90    0.12
decades_since_1950 -0.36    0.78
margin             -1.14    0.68
---
n = 153, k = 5
residual sd = 11.46, R-Squared = 0.31

```

Now we're getting somewhere! The estimate is 5–10 years, and it's 2 standard errors from zero. We can juice it up a bit more by removing `decades_since_1950` (a reasonable choice to remove) and `age`. We should really keep `age` in the model, no question about it, not including `age` in a remaining-length-of-life model is about as bad as not including smoking in a cancer model, but let's remove it just to see what happens:

```

lm(formula = more_years ~ won + margin, data = data, subset = subset4)
      coef.est coef.se
(Intercept) 22.70    2.52
won          11.92    4.31
margin       -2.00    0.77
---
n = 153, k = 3
residual sd = 13.34, R-Squared = 0.05

```

Finally! An estimated effect of more than 10 years of life, *and* it's 2.8 standard errors from zero.

So now we can reverse engineer, starting with the regression discontinuity model from the paper and keeping only the first race for each candidate, which yields:

Method	Coef.	Std. Err.	z	P> z	[ 95% C.I. ]
Conventional	1482.670	1021.784	1.451	0.147	[-519.991 , 3485.330]
Robust	-	-	1.516	0.130	[-525.718 , 4115.893]

If we then include `age` as a predictor in the `rdrobust` call, we get:

Method	Coef.	Std. Err.	z	P> z	[ 95% C.I. ]
Conventional	819.480	708.968	1.156	0.248	[-570.072 , 2209.032]
Robust	-	-	1.051	0.293	[-740.342 , 2453.552]

The robust setting in this package doesn't matter much in this example—but the analysis is sensitive to the bandwidth (yes, the bandwidth is estimated from the data, but that just tells us it can be noisy; the fact that something is calculated automatically using some theory and a computer program doesn't mean it's correct in any particular example) and to the decision of how to handle candidates with multiple elections in the dataset, and to the modeling of the `age` predictor.

## 4. Discussion

### 4.1. No smoking gun

It's easy to criticize research when the forking paths are all out in the open (as with the Bem, 2011, study of ESP) or when statistics show that your sample size is too low to detect anything by a factor of 100 (as in the Kanazawa, 2007, study of beauty and sex ratio; see Gelman and Weakliem, 2009) or when there are obvious forking paths and failed replications (as in many notorious studies in social psychology) or when almost all the data have been excluded from the analysis (as in the

study by Kim, Zhang, and Zhong, 2021, of unionization and stock prices; see Gelman, 2019) or when there's flat-out research misconduct.

This example discussed here is a bit different. It's a clean analysis with clean data. The data are even publicly available (which allowed me to make the above graphs), and the researchers responded directly to my concerns. But, remember, honesty and transparency are not enough. If you do a study of an effect that is small and highly variable (which this one is: to the extent that winning or losing can have large effects on your lifespan, the effect will surely vary a lot from person to person), you've set yourself up for scientific failure: you're working with noise.

I'm not happy about this, but that's just how quantitative science works. So let me emphasize again that a study can be fatally flawed just by being underpowered, even if there's no other obvious flaw in the study.

Or, to put it another way, there's an attitude that causal identification + statistical significance = discovery, or that causal identification + robust statistical significance = discovery. But that attitude is mistaken. Even if you're an honest and well-meaning researcher who has followed principles of open science.

#### **4.2. Awkwardness of this exercise**

This was all a lot of work. I did it because it's worth doing some work now and then—but don't forget the larger point, which is that we were suspicious from the start because of the implausibly large effect sizes, and we knew ahead of time not to assume that causal identification + statistical significance + robustness tests = discovery, because we've been burned on that combination many times before.

To put it another way, don't think that you should give the claims of a published or preprinted social science study the benefit of the doubt, just because someone like me didn't bother to go to all the trouble to explain its problems. The point here is not to pick on the authors of this particular study—not at all. They're doing what they've been trained to do.

Again, nothing special about regression discontinuity here. All the same concerns would arise with any observational study, whether it uses matching, difference in differences, nonparametric modeling, instrumental variables, synthetic controls, or just plain regression. The same problems arise in experiments too, what with issues of missing data and extrapolation. Indeed, sometimes it's worse with designed experiments, because of the illusion that they give clean causal identification. With all these kinds of study, if the underlying effect size is small and your measurements are noisy, you're drawing dead (as they say in poker). Again, this can happen in any study, ranging from a clean randomized experiment at one extreme, to a pure observational study on the other.

#### **4.3. Potential rebuttal**

As noted, authors and different groups of readers collaborate implicitly and asynchronously, and so it can help to move forward by anticipating objections to one's work. To this end, I constructed the following response by a hypothetical skeptic of my skepticism:

Hey, destructive statistician. Stop with your nihilism! The authors of this peer-reviewed paper did a very standard, reasonable analysis and found a large, statistically significant, and robust effect. All you did was perform a bunch of bad analyses. Even your bad analyses uniformly found positive effects. They just weren't statistically significant, but that's just cos you threw away data and used crude methods. It's well known that if you throw away data and use inefficient statistical analysis, that your estimates become more noisy, your standard errors become larger, and your power has gone down. In

summary, you have maligned a competent, high-quality paper—one that passed peer review—by the simple expedient of re-running their analysis with fewer data points and a noisier statistical method and then concluding that their results were not statistically significant.

The above response sounds kind of reasonable—indeed, readers might agree with it!—and it has the form of a logical argument, but I think it’s wrong, for several reasons.

First, regarding the point that the coefficient shows up positive in all my analyses, hence making this nothing more than a dispute over statistical significance: This where it’s helpful to have a Bayesian perspective—or, if you’d prefer, a design-based perspective. If the true effect is in the range of 1 year of life (in either direction), and this is studied with data and analysis that yield an estimate with standard errors of 2 years or more, then the signal will be dominated by noise, and a positive coefficient can be explained as just the summation of some highly variable numbers that happened to end up in one direction or another.

Second, regarding the idea that I’ve replaced their state-of-the-art regression discontinuity analysis with various low-power linear regressions: Actually, no. The robust adaptive discontinuity regression does not have higher power or more statistical efficiency than the linear model. These are just different models. Indeed, the decision of the authors of the original paper to not include age as a predictor *lowers* the power of the analysis. Regarding sample size: I took out the duplicate elections featuring the same candidate because it’s inappropriate to consider these as independent data points; indeed, including them in that way can just give you artificially low standard errors. (I think it would be easy enough to show this using a bootstrap analysis or something like that.) The real point, though, is you can’t tell much about the power of a study by looking at the  $z$ -score or statistical significance of a coefficient estimate.

#### 4.4. Reactions of the original authors

It is good to anticipate criticism but better to actually receive it. In this case, the authors of the paper under discussion had contacted me directly, and after I posted my reactions online, one of them responded (Larsen, 2020), and I posted a response to their response (Gelman, 2020). We had some back-and-forth about the details of how I attempted to reproduce the published analysis, but the most important sticking point was the plausibility of the effect size. In an email correspondence, Larsen wrote:

If the “true” effect is, say, .5 years, we would need a lot more data to detect that. In other words, I believe your point 4 is spot on. That’s indisputable and simply a limitation of the paper. I am not sure I would say that the paper is ‘fatally flawed’ for that reason, but I get the point. Are you saying that we might as well have detected a negative effect of 5-10 years? (I mean, if we’re simply playing around with noise here.) ... We had a lot of discussions about the effect size when writing up the paper and, to be honest, I don’t know what’s a realistic/common sense effect size to expect here. When you work on a paper long enough, your results become ‘common sense’ (or to paraphrase Duncan Watts ... ‘everything is obvious once you know the answer’).

I replied:

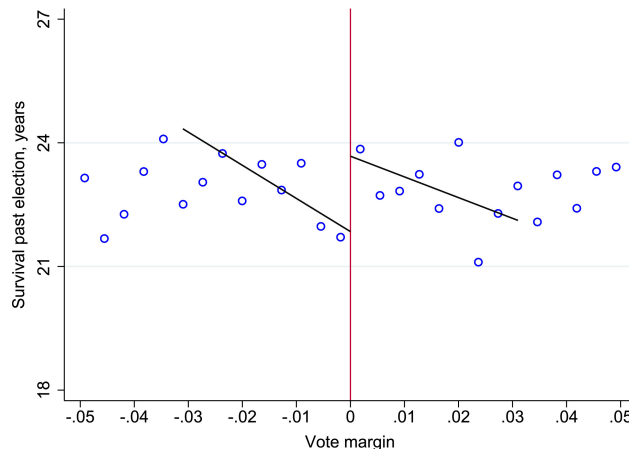
I suspect any true effect would be much less than 5 years, so that a pattern in the data that shows up in the regression is basically just noise. I think that in a replication study of other data with the same sample size, for example from some state legislatures or whatever, it would be possible to get an estimate of  $-5$  yrs or  $-10$  yrs. I expect that

a fully preregistered study would not yield a statistically significant result, but your study was not preregistered (nor are most of mine! I'm not holding my own research up as a model here) ... Is an effect size of 5 or 10 years reasonable? I agree that we can't know, but it really doesn't make sense to me, especially given that candidates typically face many elections in their lifetimes. The point of the type M errors is that, with such a small sample and such noisy data, anything statistically significant estimate you will find will *have* to be huge. So, in that sense, I don't see your estimate of 5–10 yrs as providing any useful information regarding the magnitude of the effect: Even it turns out that there is some average effect, even as large as 1 or 2 years (which I would really doubt), the observed estimate of 5–10 yrs is an artifact of the small sample and high variability.

Larsen added:

Also, a good friend of mine informed me (when he saw the study) that we're not the first to show an interest in this subject. Apparently, a couple of economists published a study on the topic in 2019 as well. They find a positive effect but a lot smaller, i.e. 3 years when looking at elections post-1908. Interestingly, they find that governors running in elections post-1908 (the sample most similar to ours) live  $\sim 6$  years longer (Panel B, Column 1–2, Table 3). I am not linking to this study to make a point about the results (again, I believe you have a good point on the fact that we would need more statistical power to estimate a smaller effect), but it might be of interest to cite that study as well to make a broader point about inference and robustness?

I took a look at the study from 2019 (Borgschulte and Vogler, 2019), and here is their plot:



This is strikingly similar to the main result of Barfort et al. (see the first graph reproduced in this article), which demonstrates the way in which regression discontinuity analyses have sufficient researcher degrees of freedom to allow statistically significant coefficients to be extracted from data that are consistent with pure noise, as discussed by Gelman (2017). And this returns us to my skepticism about a mechanism by which losing elections would reduce candidates lifespans by 5 or more years on average.

#### 4.5. Reading and using research articles as a form of asynchronous collaboration

We can assume that, when a research paper is successful, that most of its readers probably have no direct connection with the authors. And the afterlife of a published paper is not just with its

readers but also with later researchers who apply its findings or methods. As such, post-publication criticism is an important part of the feedback mechanism of science, and we can see it as a form of collaboration occurring across time as well as space.

Direct, or synchronous, collaboration is characterized by frequent “handshaking”: theoretically redundant calculations or observations that allow the different members of a research team to check their answers against each other. This kind of confidence-building is more challenging in the more general setting when there is no contact between the original researchers and the replicators. The present article gives a sense of how this process can go.

## References

- Barfort, S., Klemmensen, R., and Larsen, E. G. (2021). Longevity returns to political office. *Political Science Research and Methods* **9**, 658–664.
- Bem, D. J. (2011). Feeling the future: Experimental evidence for anomalous retroactive influences on cognition and affect. *Journal of Personality and Social Psychology* **100**, 407–425.
- Borgschulte, M., and Vogler, J. (2019). Run for your life? The effect of close elections on the life expectancy of politicians. *Journal of Economic Behavior & Organization* **167**, 18–32.
- Gelman, A. (2017). Forking paths plus lack of theory = No reason to believe any of this. <https://statmodeling.stat.columbia.edu/2017/12/29/forking-paths-plus-lack-theory-no-reason-believe/>
- Gelman, A. (2018). The failure of null hypothesis significance testing when studying incremental changes, and what to do about it. *Personality and Social Psychology Bulletin* **44**, 16–23.
- Gelman, A. (2019). Another regression discontinuity disaster and what can we learn from it. <https://statmodeling.stat.columbia.edu/2019/06/25/another-regression-discontinuity-disaster-and-what-can-we-learn-from-it/>
- Gelman, A. (2020). Comment. <https://statmodeling.stat.columbia.edu/2020/07/02/no-i-dont-believe-that-claim-based-on-regression-discontinuity-analysis-that/#comment-1373309>
- Gelman, A., and Imbens, G. (2017). Why high-order polynomials should not be used in regression discontinuity designs. *Journal of Business and Economic Statistics* **37**, 225–456.
- Gelman, A., and Loken, E. (2014). The statistical crisis in science. *American Scientist* **102**, 460–465.
- Gelman, A., and Weakliem, D. (2009). Of beauty, sex, and power: Statistical challenges in estimating small effects. *American Scientist* **97**, 310–316.
- Gelman, A., and Zelizer, A. (2015). Evidence on the deleterious impact of sustained use of polynomial regression on causal inference. *Research and Politics* **2**, 1–7.
- Kanazawa, S. (2007). Beautiful parents have more daughters: A further implication of the generalized Trivers-Willard hypothesis. *Journal of Theoretical Biology* **244**, 133–140.
- Kim, J. B., Zhang, E. X., and Zhong, K. (2021). Does unionization affect the managershareholder conflict? Evidence from firm-specific stock price crash risk. *Journal of Corporate Finance* **69**, 101991.
- Nosek, B. A., Spies, J. R., and Motyl, M. (2012). Scientific utopia II. Restructuring incentives and practices to promote truth over publishability. *Perspectives on Psychological Science* **7**, 615–631.

- Larsen, E. G., (2020). A response to Andrew Gelman. <https://erikgahner.dk/2020/a-response-to-andrew-gelman/>
- Simmons, J., Nelson, L., and Simonsohn, U. (2011). False-positive psychology: Undisclosed flexibility in data collection and analysis allow presenting anything as significant. *Psychological Science* **22**, 1359–1366.