### A Quantitative Method for Substantive Robustness Assessment

Forthcoming in Political Science Research and Methods

Justin Esarey\* and Nathan Danneman<sup>†</sup>
April 10, 2014

#### Abstract

Empirical political science is not simply about reporting evidence; it is also about coming to conclusions on the basis of that evidence and acting on those conclusions. But whether a result is substantively significant—strong and certain enough to justify acting upon the belief that the null hypothesis is false—is difficult to objectively pin down, in part because different researchers have different standards for interpreting evidence. Instead, we advocate judging results according to their "substantive robustness," the degree to which a community with heterogeneous standards for interpreting evidence would agree that the result is substantively significant. We illustrate how this can be done using Bayesian statistical decision techniques. Judging results in this way yields a tangible benefit: false positives are reduced without decreasing the power of the test, decreasing the error rate in published results.

# Introduction: statistical inference and rational choice under uncertainty

A long and and cross-disciplinary literature stresses the importance of assessing the substantive significance of empirical results (Achen, 1982; Hunter, 1997; McCloskey, 1998; Gill, 1999; Lunt, 2004; Ziliak and McCloskey, 2004, 2008; Miller, 2008; Siegfried, 2010); we say that a result is "substantively significant" when it is justifies acting on the belief that the

<sup>\*</sup>Assistant Professor, Department of Political Science, Rice University. Corresponding author (justin@justinesarey.com).

<sup>&</sup>lt;sup>†</sup>Department of Political Science, Emory University. E-mail: (ndannem@emory.edu).

null hypothesis is false.<sup>1</sup> Naturally, substantive significance of this sort is partially a matter of individual judgment. These standards include how averse a researcher is to drawing a mistaken conclusion and how large a relationship must be in order to be scientifically or politically meaningful. Thus, attempts to use decision theory (Wald, 1950; DeGroot, 2004 {1970}; Pratt, Raiffa and Schlaifer, 1996; Manski, 2007, ch. 12) to establish criteria for substantative significance are likely to degenerate into disagreement over the utility function chosen to encapsulate such standards.

In this paper, we argue that it is more helpful to assess a result's substantive robustness, the degree to which a community with heterogeneous standards for interpreting evidence would agree that the result is substantively significant, rather than whether it meets any individual standard. This focuses attention away from contention about which specific standards are appropriate and onto the breadth of evaluation standards that can be satisfied by a particular piece of evidence. The idea is to enable a researcher to objectively demonstrate whether his/her results should be regarded as substantively significant by a scientific community, even if there is significant disagreement in that community over what a substantively significant result looks like. This approach is compatible with any number of sufficiently flexible utility functions, but we offer one that we believe is well-suited to the task.

We also offer evidence that requiring results to be substantively robust as well as statistically significant could improve the quality of hypothesis tests in the discpline. Our Monte Carlo studies reveal that:

- 1. Results that are both statistically significant and substantively robust are less likely to be false positives<sup>2</sup> than those that are merely statistically significant, with no loss in power.<sup>3</sup>
- 2. Statistically significant but substantively meaningless results are extremely unlikely

<sup>&</sup>lt;sup>1</sup>We thank an anonymous reviewer for suggesting this phraseology.

<sup>&</sup>lt;sup>2</sup>By "false positive," we mean incorrectly rejecting a one-sided null hypothesis when that hypothesis is true

 $<sup>^{3}</sup>$ By "power," we mean the probability of rejecting a one-sided null hypothesis when that hypothesis is false.

when a scientifically important relationship actually exists; observing this pattern allows us to confidently conclude that the result is a false positive.

In short, substantive robustness tests may enable researchers to discriminate between genuine results and statistical anomalies.

As has been noted in the past (Ziliak and McCloskey, 2004), the statistical significance of a result can be markedly more or less robust than its substantive significance. We illustrate this with a re-analysis of two recently published results. In a recent piece by Clinton (2006), we find that some effects that are statistically significant are not substantively robust to a reasonable range of standards for the interpretation of evidence. We also re-analyze data from Clark and Golder (2006), confirming the substantive robustness of the authors' results while demonstrating the application of our technique in situations where multiple variables determine the effect of interest (such as in models with interacted variables).

### Drawing conclusions from evidence

To begin, we present an example to help the reader crystallize the dimensions of judgment on which researchers may differ when making judgments of substantive significance. We think that this is easiest to do by way of an analogy between an intuitively compelling but evidence-based decision, and a more ordinary social scientific decision that a researcher might make in day-to-day work.

Figure 1 illustrates this analogy. Let us assume that researchers develop an interesting new drug: it can reliably produce any level of weight loss desired by the user. But researchers suspect that a fatal cancer might be a side effect of the drug. Panel 1a of Figure 1 depicts a hypothetical posterior belief distribution of the change in cancer probability based on a series of (imagined) studies; we denote this change in cancer probability r. The distribution is centered over a 25% increase in a subject's chance of developing cancer, but the 95% confidence interval is [-48%, 98%] corresponding to a one-tailed p-value of about 0.25. That

is, the relationship between this drug and cancer is statistically insignificant.

Would you take this drug?

A rational chooser<sup>4</sup> would consider the potential consequences of taking the drug vs. not taking the drug under each of the possible states of the world—each of the possible values of  $\partial \Pr(\text{cancer})/\partial \text{take drug} = r$  under the posterior—and then determine whether taking the drug has a positive expected utility:

$$E[u(\text{take drug}) - u(\text{don't take drug}) \mid \text{data}] = \int [u(\text{take drug}|r) - u(\text{don't take drug}|r)] f(r|\text{data}) dr = \int u(\text{take drug}|r) f(r|\text{data}) dr - \int u(\text{don't take drug}|r) f(r|\text{data}) dr = E[u(\text{take drug}) \mid \text{data}] - E[u(\text{don't take drug}) \mid \text{data}]$$
(1)

The decision to choose to take the drug when the expected utility (in equation 3) is greater than zero, and to not take it otherwise, is the *Bayes decision rule*:

$$d(x)$$
 = take drug if  $E[u(\text{take drug})|x] > E[u(\text{don't take drug}|x)]$   
don't take drug otherwise

where we let x signify the data. The rule associates each state of the world f(r|data) with an optimal Bayes action that maximizes expected benefit (French and Insua, 2000, p. 148). Uncertainty about the actual effect of the drug plays a role in the decision<sup>5</sup> through f(r|data). This Bayes decision rule d(x) is optimal in that it minimizes the Bayes risk, the expected

$$f(r|\text{data}) = \frac{f(\text{data}|r)f(r)}{\int f(\text{data}|r)f(r)dr}$$

The analyst must specify a prior, f(r), before determining whether these conditions are met; we will discuss prior specification in a succeeding subsection.

<sup>&</sup>lt;sup>4</sup>See DeGroot (2004 {1970}, Chapter 7), French and Insua (2000, Chapter 6), and Pratt, Raiffa and Schlaifer (1996, Chapters 3 and 4) for general overviews of Bayesian statistical decision theory.

<sup>&</sup>lt;sup>5</sup>The "Bayesian" in Bayesian statistical decision theory comes from the fact that these conditions rely on the posterior probability of r given the data, f(r|data), which is determined by Bayes' rule:

loss in utility  $l(\bullet) = -u(\bullet)$  from a decision rule averaging over both sampling variation in the data and the prior distribution of r:

$$BR(d) = \int \left[ \int l(d(x)|r) f(x|r) dx \right] f(r) dr$$

Summarily, it is rational to commit to following the Bayes decision rule because it maximizes utility (equivalent to minimizing loss) given all sources of uncertainty in our empirical evidence.<sup>6</sup>

While there may be disagreements over the quality of the evidence, we set them aside to focus on differences in interpretation of the same evidence. These differences in judgment presumably come through the utility function u. For example, some people may be more loss averse<sup>7</sup> than others; they would presumably be less willing to gamble on the possibility of cancer in order to lose weight. The difference on this dimension is in how people weight the relative value of forsaking a gain if the drug is harmless compared to the value of avoiding a loss if the drug causes cancer. People may also differ in their absolute valuation of a cancer outcome: a small but completely certain increase in cancer incidence (say, 0.5%) might be

$$BR(d) = \int \left[ \int l(d(x)|r) f(x|r) dx \right] f(r) dr$$

Noting that the central term  $R(x,\beta) = \int l(d(x)|r) f(x|r) dx$  is often referred to as the (frequentist) risk, the expected loss integrating over sampling variation given a fixed state of the world r. Reverse the order of integration without loss of generality:

$$= \int \left[ \int l(d(x)|r) f(x|r) f(r) dr \right] dx \tag{2}$$

This implies that a decision rule which minimizes the expected loss for every particular data set,  $\int l(d(x)|r) f(x|r) f(r) dr$ , will also minimize the Bayes risk. But this is precisely what the Bayes decision rule prescribes; it minimizes:

$$\int l(d(x)|r)f(r|x)dr = \frac{\int l(d(x)|r)f(x|r)f(r)dr}{\int f(x|r)f(r)dr}$$

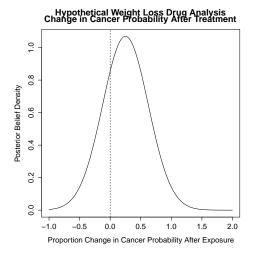
where the numerator is the central term of equation 2 and the denominator is a constant.

<sup>&</sup>lt;sup>6</sup>A proof of this proposition is given by French and Insua (2000, pp. 164-165). Begin with the expression of Bayes risk:

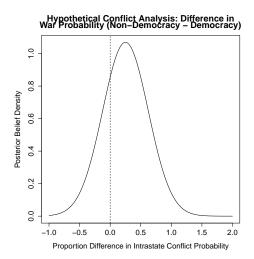
<sup>&</sup>lt;sup>7</sup>Here *loss aversion* refers to overvaluing reductions in utility relative to gains when computing expected values; this is not directly related to the similarly-named *loss function* (i.e., the negative of the utility function in equation 3).

Figure 1: Hypothetical Evidence-Based Decisions

(a) Would You Take This Weight Loss Drug?



(b) Would You Believe that Non-Democracies Have More Civil Wars than Democracies?



negligible to some compared to the benefit of weight loss, while others may consider this same incidence extremely meaningful from a personal or public health perspective.

Now, consider another decision: should we act on the conclusion that democratic states are less likely to suffer a civil war compared to non-democracies, or treat democracies at least equally likely to suffer a civil war? Based on (hypothetical) evidence shown in Panel 1b of Figure 1, the posterior distribution of the difference in intrastate conflict probability (which we denote as  $\beta$ ) is the same as the distribution of the difference in cancer probability from Panel 1a. As before, a rational chooser should determine whether there is greater scientific (or policy) benefit to acting on the conclusion that democracies are less susceptible than democracy to civil war, or in treating democracies as at least equally susceptible, based on the available evidence:

$$E[u(\text{treat as unequal}) - u(\text{treat as equal})] = \int [u(\text{treat as unequal}|\beta) - u(\text{treat as equal}|\beta)] f(\beta|\text{data})]d\beta \quad (3)$$

More generally, we could write this decision as:

$$E[u(\text{act on alternative}) - u(\text{act on null})] = \int [u(\text{act on alternative}|\beta) - u(\text{act on null}|\beta)] f(\beta|\text{data})]d\beta \quad (4)$$

where the alternative hypothesis is that democracies are less prone to civil war and the (default) null hypothesis is that they are at least equally prone to civil war. Of course, we would not expect  $u(\text{act on alternative}|\beta) - u(\text{act on null}|\beta)$  to be the same for this decision as for the decision to take the weight loss drug. But the considerations involved in shaping the judgment are similar: (1) how much does the scientist weight the consequences of falsely acting on the alternative against the consequences of falsely acting on the null, and (2) how large must democracies' advantage in civil war susceptibility be before it merits scientific and political attention?

Rather than argue for a particular answer to these questions, our goal is to determine whether researchers with a variety of different standards would make the same decision according to equation 4. We must still choose a utility function, but ideally it will be one that can fairly represent a variety of standards for evaluating evidence. In particular, it should be able to represent a continuum of relative valuations for false positives (mistakenly acting on the alternative hypothesis) and false negatives (mistakenly acting on the null hypothesis).

### Quantifying the decision: utility, loss aversion, and decision rule

The functional form of  $u(\text{act on alternative}|\beta) - u(\text{act on null}|\beta)$  is the encapsulation of scientific judgment in this framework, and there is no universally correct choice. Recall that our objective is to find whether different standards of judgment would lead to the same conclusion about the substantive significance of an empirical result. Thus, the precise choice of functional form is less important compared to its ability to approximate a diversity

of preferences about (a) the tradeoff between false positives and false negatives, and (b) the minimum effect size required for substantive significance. We suggest a function that is simple, yet flexible enough to adapt to varying decision contexts and the preferences of different researchers; we do not claim that it is unique or superior to all possibilities, but it meets our stated criteria well and provides some inferential benefits that we explore in a later section.

For the decision of whether to act on the conclusion that the relationship between two variables (such as democracy and intrastate conflict) is positive and large enough to be scientifically or politically important, we focus on a loss averse utility function:

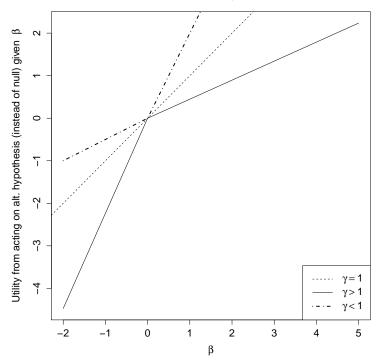
$$u(\text{act on alternative}|\beta, \gamma, c) - u(\text{act on null}|\beta, \gamma, c) = \gamma^{-\text{sign}(\beta - c)} [\beta - c]$$
 (5)

Here, the alternative is the conclusion that  $\beta \geq c$  and the null is the conclusion that  $\beta < c.^8$  This utility function is depicted in Figure 2 for c=0. In terms of our previous democracy-conflict example, this is equivalent to saying that a non-democracy's chance of suffering a civil conflict is at least c more than a democracy's. Utility is gained in proportion to distance from the threshold for substantive significance c: with all else held equal (i.e., in a single decision context), bigger relationships are more substantively important than small ones and the smallest positive relationship worth paying attention to is of size c. That is, the consequences of accepting this conclusion are linearly related to the size of the difference.  $\gamma$  is a parameter that describes the valuation ratio for correct and incorrect decisions, or the degree of loss aversion that the researcher exhibits;  $\gamma$  is the degree of kink in the function in Figure 2. Downward kinks indicate that correctly acting on the alternative is less important than incorrectly doing so (a false negative is less damaging than a false positive), and correspond to  $\gamma > 1$ . Upward kinks indicate that correctly acting on the alternative is more important than incorrectly doing so (that is, a false positive is less damaging than a false negative) and

<sup>&</sup>lt;sup>8</sup>The utility from accepting an alternative hypothesis of a negative relationship uses the same functions below, but replaces  $\beta$  with  $-\beta$  and c with -c so that negative  $\beta$  values yield positive utility and so that the substantive threshold for significance is less than zero.

Figure 2: Loss Averse Utility Functions

#### **Loss Averse Utility Function**



correspond to  $\gamma \in (0, 1)$ .

The Bayes decision rule under this utility framework is to act on the alternative hypothesis  $(\beta \geq c)$  whenever:

$$\int \gamma^{-\operatorname{sign}(\beta-c)} \left[\beta - c\right] f(\beta|\operatorname{data}) d\beta > 0 \tag{6}$$

and to act on the null hypothesis otherwise.  $^9$  There are two unknowns,  $\gamma$  and c, that

$$R(d,\beta) = \int [\gamma^{-1} [\beta - c] * I(\beta > c) * \Pr(\Delta(x) < 0|\beta) - \gamma [\beta - c] * I(\beta \le c) * \Pr(\Delta(x) > 0|\beta)] dx$$

Where we set u (null| $\beta$ <c) = u (alternative| $\beta \ge c$ ) = 0, u(null| $\beta \ge c$ ) =  $-\gamma^{-1}$  [ $\beta - c$ ], and u(alternative| $\beta < c$ ) =  $\gamma$  [ $\beta - c$ ]; recall that  $l(\bullet) = -u(\bullet)$ . The Bayes risk further integrates over the distribution of  $\beta$  using the decision rule that forms the basis for R:

$$BR(d) = \int R(x,\beta)f(\beta)d\beta$$

Given that the decision rule is (by construction) a Bayes decision with respect to the prior distribution  $f(\beta)$ , it minimizes the Bayes risk and is also f-admissible (French and Insua, 2000, Proposition 30 on p. 166; see

<sup>&</sup>lt;sup>9</sup>Letting x stand for the data and  $\Delta(x) = \int \gamma^{-\operatorname{sign}(\beta-c)} [\beta-c] f(\beta|x) d\beta$  from equation (6), the frequentist risk  $R(x,\beta) = \int l(d(x)|\beta) f(x|\beta) dx$  of this rule is

correspond to the two variable dimensions of scientific judgment that we laid out at the start. By determining the  $c^*$  that solves:

$$\int \gamma^{-\operatorname{sign}(\beta - c^*)} \left[ \beta - c^* \right] f(\beta | \operatorname{data}) d\beta = 0$$
 (7)

we can sketch a  $\gamma/c^*$  curve over which the researcher is indifferent to acting on the alternative hypothesis given  $\gamma$  and  $c^*$ .

### Computing $c^*$ and choice of a prior

also p. 144)

Computing the  $\gamma/c^*$  curve is a process that the vast majority of mainstream statistical software packages are already capable of through add-on packages. We have developed software for R and Stata that, when given a  $\gamma$ , will compute a  $c^*$  immediately after a linear regression or other generalized linear model. The software will also determine  $c^*$  when given manually entered information from a published table (without the accompanying data set).

The software package assumes a truncated uniform prior distribution  $f(\beta)$  on the interval defined by  $\hat{\beta} \pm 8\sigma$ . This structure assumes minimalistic knowledge of the underlying parameters before examining a data set. It also allows for a Bayesian interpretation of frequentist results: for the classical linear regression model,  $f(\beta|\text{data})$  takes a multivariate t distribution with n-k degrees of freedom (Gelman et al., 2003, pp. 355-357). The software therefore uses a t distribution for the posterior when computing  $c^*$  in these cases. For generalized linear models with this prior estimated via maximum likelihood,  $f(\beta|\text{data})$  is asymptotically normal as  $n \to \infty$  (Gelfand and Ghosh, 2000, pp. 4-8), and therefore the software uses the normal distribution in these instances.

Computing  $c^*$  is a matter of finding the roots of equation 7 (or the equivalent for an alternative utility function). This class of problem is already solved via iterative maximization algorithms, such as Newton-Raphson, 10 but we are required to compute an integral at

<sup>&</sup>lt;sup>10</sup>Many maximization algorithms involve finding a root of  $\frac{df(x)}{dx}$  to find a maximum of f(x). The root-

each iteration because we are calculating an expected value rather than a maximum or modal value. For results with a t-distributed or normally distributed posterior, we can use standard quadrature approaches to quickly compute the integral for each step of the maximization process.

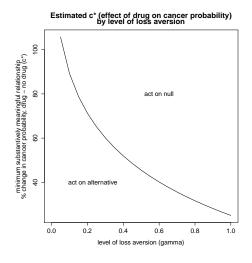
Some analysts may wish to impose stronger prior beliefs, which will change the distribution of  $f(\beta|\text{data})$ . If the prior is conjugate with the posterior, the analyst must replace the multivariate t or normal distributions with the appropriate posterior distribution before calculating  $c^*$ . For more exotic posteriors, where  $f(\beta|\text{data})$  is difficult to analytically express, an approximation can be computed via Markov Chain Monte Carlo methods, stored, and then used in the integration of equation 7.

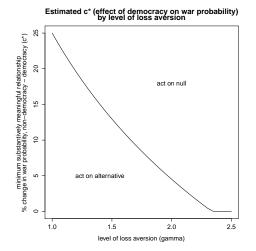
## What does a substantively robust (or non-robust) finding look like?

Figure 3b sketches the  $\gamma/c^*$  curve for our democracy and war example using equation 7. The relevant question is: can we conclude that non-democracies are susceptible to greater war risk on the basis of the evidence from Figure 1b? Our answer is that a researcher would accept an increased civil war probability under the point on the curve corresponding to his/her  $\gamma$ . Equivalently, a researcher would decide how large a relationship had to be before it was scientifically or politically meaningful and how loss-averse s/he is, then find the appropriate point in Figure 3b and make the decision indicated. For example, a person who weighted false positives and false negatives equally (with  $\gamma = 1$ ) would be willing to act on the belief that non-democracies are  $\leq 25\%$  less likely to experience civil war than comparable non-democracies. But a person who valued false positives four times as much as false negatives  $(\gamma = \sqrt{4} = 2)$  would only be willing to act on a  $\leq 5\%$  difference in civil war risk, and a person who valued false positives even more highly would always act on the null (of a nonexistent finding capability of these algorithms is easily adapted to non-maximization root solutions.

Figure 3: Substantive Robustness of Inference from Figure 1

- (a) Drugs and Cancer Probability (Figure 1a)
- (b) Democracy and War Probability (Figure 1b)





or negative relationship). In short, the association between democracy and civil war risk (in this hypothetical example) is not robust to reasonable ranges of variation in standards for the interpretation of evidence. It is not necessarily important for a researcher to know his or her own  $\gamma$ , but simply to acknowledge that a community of people with reasonable variation in how much they valued false positives more highly than false negatives would not agree upon the substantive significance of the relationship between democracy and civil war risk.

By contrast, the conclusion that we should not take the weight-loss drug (for fear of cancer risk, as shown in Figure 1a) is quite robust. Figure 3a sketches the  $\gamma/c^*$  curve for this example; the key difference, compared to the democracy/war example, is that  $\gamma \in [0, 1]$  would be appropriate for a person who was much less concerned about mistakenly acting on the alternative (accidentally concluding that the drug causes cancer when it does not, and therefore not taking the drug when it would be beneficial) compared to mistakenly acting on the null (mistakenly believing that the drug does not cause cancer when it actually does, and therefore taking the drug when it increases one's likelihood of developing cancer). For this reason, even wide variation in standards for the interpretation of evidence produces the conclusion that there is a large and substantively meaningful link between the drug and

cancer risk. It is not so important for any particular person to locate their personal value of  $\gamma$ , but simply to acknowledge that it lies somewhere between 0 and 1 (i.e., that they value false negatives more highly than false positives in some measure) and that people with diverse preferences in this domain would avoid taking the drug.

We think these examples illustrate the value of focusing on the robustness of substantive significance to different standards rather than on any particular judgment of substantive significance. Figure 3 makes clear that comparatively small differences in loss aversion (that is, aversion to falsely rejecting the null) translate to large differences in the assessment of substantive significance. Aggressive researchers who are relatively indifferent to false positives  $(\gamma \approx 1)$  would conclude (based on our hypothetical evidence) that there is a substantively significant association between democracy and civil war probability as long as they thought that a 25% increase in war probability was meaningful. More cautious researchers who weight false positives as more damaging ( $\gamma \gtrsim 2.5$ ) would never draw this conclusion. These are rather wide differences in judgment that relate to incommensurable differences in scientific viewpoint. Consequently, we think that arguing for a universally acceptable  $\gamma$ , let alone a universally acceptable form for u(act on alternative) - u(act on null), is probably hopeless in the context of academic research. By focusing on the breadth of preferences for which a result is substantively significant rather than any particular preference, we sidestep the knotty question of "whose judgment is correct?" in favor of the more readily answerable "would this evidence satisfy most researchers?"

# Substantive robustness assessment can improve the quality of published results

When it comes to statistical hypothesis testing, we determine that it is advantageous to require a result to be both substantively robust and statistically significant in order to be scientifically notable, compared to statistical significance alone. Aside from the theoretical

and descriptive value of assessing the substantive robustness of results, requiring a statistical relationship to be statistically significant and substantively robust before rejecting the null hypothesis (of no relationship) results in a more powerful test (a greater probability of rejecting false null hypotheses) at any given size (probability of rejecting correct null hypotheses). This advantage accrues to the combined procedure because results that are statistically significant but substantively fragile frequently occur when the null hypothesis is true, but rarely occur when the alternative is true. The upshot is that a researcher using our procedure will draw fewer mistaken conclusions than a researcher just using statistical significance testing.

To demonstrate, we conduct a simulation to calculate the rate at which the null hypothesis is rejected using two different procedures:

- 1. statistical significance: reject if two-tailed  $p \leq \alpha$ .
- 2. combined statistical significance and substantive robustness: reject if
  - (a) result is statistically significant (two-tailed  $p \leq \alpha$ ), and
  - (b)  $c^* > 0 \text{ for } \gamma = 2.$

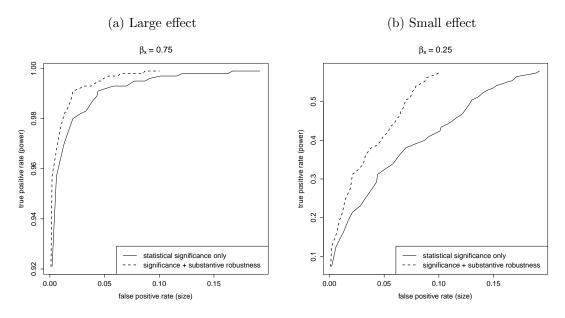
For the second procedure, we choose a relatively small  $\gamma=2$  in order to ensure a very minimal bar for substantive robustness. We also set the critical threshold for substantive acceptability to  $c^*>0$ , the smallest (and easiest-to-pass) threshold for a test of a positive relationship, for the same reason. In short, a result will be judged substantively robust in this simulation if someone with  $\gamma \leq 2$  would be willing to act on the alternative hypothesis, even for a very small relationship size.

Our simulation uses a linear DGP:

$$y = 2 + \beta_x x + u$$

In all simulations,  $u \sim \Phi(\mu = 0, \sigma = 0.5)$ . We use three values of  $\beta_x$ :  $\beta_x = 0.75$  (a "large effect" more easily detected by statistical tests),  $\beta_x = 0.25$  (a "small effect" harder to

Figure 4: Size/Power Analysis for Statistical Significance Testing, With and Without Substantive Robustness Checking



distinguish from noise), and  $\beta_x = 0$  (the null hypothesis). For each value of  $\beta_x$ , we conduct 1000 simulations with data sets of size n = 100 each. We then repeat these simulations for values of  $\alpha$  (the critical value for the statistical significance test) between 0.005 and 0.2 to get a range of false positive rates. The results are depicted in Figure 4.

Figure 4a plots the proportion of the time that a test procedure rejected the null when  $\beta_x = 0$  on the x-axis, and the proportion of the time that the test rejected the null when  $\beta_x = 0.75$  on the y-axis; Figure 4b plots the same relationship, but with the null rejection rates for  $\beta_x = 0.25$  on the y-axis. Ideal performance would mean a completely  $\Gamma$ -shaped curve—an estimator that could achieve a 100% true positive rate at a 0% false positive rate. Better testing procedures are those which have larger true positive rates for every false positive rate: they can correctly reject false null hypotheses without accidentally rejecting true null hypotheses. As the figure shows, combining statistical significance with a very minimal evaluation of substantive acceptability—a mildly loss-averse researcher ( $\gamma = 2$ ) must be willing to accept an effect of any magnitude—improves the power of the test at every size. The performance gap varies according to the nature of the DGP and the  $\alpha$  value

of the statistical significance test, but combined testing improves power by as much as 10 percentage points in our simulations (for a false positive rate  $\approx 10\%$  in Figure 4b).

We also find support for a very useful rule of thumb: if a result is statistically significant but not substantively meaningful ( $c^* = 0$  when  $\gamma = 2$ ), in all likelihood the null hypothesis is the correct one. The relevant Bayesian formula for beliefs in this situation is:

$$\begin{aligned} \Pr(\text{null false}|\text{stat. sig., sub. not robust}) &= \\ & \frac{\Pr(\text{result pattern}|\text{null false}) \Pr(\text{null false})}{\Pr(\text{result pattern}|\text{null false}) \Pr(\text{null false}) + \Pr(\text{result pattern}|\text{null true}) \Pr(\text{null true})} \end{aligned}$$

But in our simulations, there were no occurrences where the result was statistically significant but substantively not robust ( $c^* = 0$  when  $\gamma = 2$ ) when the null was false; all of the false negatives in the combined test were cases where the result was substantively but not statistically significant. To ensure the robustness of this finding, we repeat our simulations for  $\beta_x = 0.25$  with greater noise in the DGP ( $u \sim \Phi(\mu = 0, \sigma = 2.5)$ ), making false negatives much more likely. We still find that this pattern of results only occured 1.5 percent of the time under the null hypothesis, leading to an updated belief of:

$$\Pr(\text{null false}|\text{stat. sig., sub. not robust}) = \frac{0.015*0.5}{0.015*0.5 + 0.028*0.5} = 0.349$$

where we determine Pr(result pattern|null true) = 0.028 by repeating our simulations for  $\beta_x = 0$  with  $u \sim \Phi(0, 2.5)$ . Even under a prior probability of 50% that the null is false, the posterior probability that the null is false is only 34.9%—nowhere close to conventional significance levels in a hypothesis test. Summarily, our simulation evidence indicates that statistically significant but substantively non-robust results are almost surely ascribable to chance.

### Assessing the substantive robustness of existing research

As a technical matter, it is straightforward to assess the substantive robustness of existing research with our technique. We demonstrate using two recently published articles from prominent general interest journals in political science. These examples allow us to show how quantitatively assessing substantive robustness helps to refine our interpretation of results. They also illustrate how our techniques are applied to models with complex marginal effects, such as those from models with interaction (product) terms, where these substantively important quantities are a combination of multiple coefficient estimates and cannot be read directly off of a coefficient table. Finally, we provide documented code for each of these examples that readers can use to replicate our results or adapt for use in their own inference problems.

### Applied example: Clinton (2006)

First, we will re-examine some of the critical results from Joshua Clinton's 2006 article on representation in Congress in the *Journal of Politics* (Clinton, 2006). In this article, Clinton examines the relationship between a survey-based ideology measure of residents of American congressional districts in the year 2000 and legislator ideal points estimated on the basis of voting records. In his OLS regression analysis, Clinton finds that there is a positive relationship between the conservatism of a Republican legislator's voting record and the degree of conservatism expressed by his/her Republican constituents. The same relationship exists between a Republican legislator's record and their Democratic constituents' ideology. The conservatism of Democratic legislators' records, by contrast, is associated with the conservatism of their Republican, but *not* Democratic, constituents.

While Clinton uses both OLS regression and an errors-in-variables regression designed to correct for measurement errors, we focus on the OLS results in this replication.<sup>12</sup> We

<sup>&</sup>lt;sup>11</sup>The legislators ideal points are estimated in Clinton, Jackman, and Rivers (2004).

<sup>&</sup>lt;sup>12</sup>The c\* approach is easily applied to EIV regression, but describing this technique would distract from

Table 1: Constituent influence on "key votes" in Congress, Table 3 from Clinton (2006)

	OLS Rep.		OLS I	OLS Dem.	
	β	s.e.	β	s.e.	
Wgt. Same Party Avg. Ideology	1.614*	.4214	.1907	.3290	
Wgt. Different Party Avg. Ideology	.6997*	.3632	2.167*	.3241	
Constant	.4602*	.1182	1.145*	.0919	

Dependent variable = legislator ideology score from Clinton, Jackman and Rivers (2004). N=222 (Republicans) and 210 (Democrats).  $R^2=0.09$  (Republicans) and 0.26 (Democrats). A \* indicates statistical significance,  $\alpha=0.05$  (one-tailed).

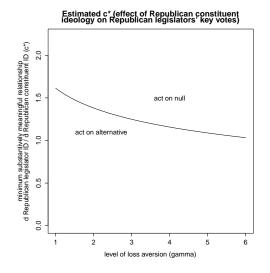
perform the same OLS regression that Clinton ran, whose results are shown in Table 1. We then calculate a substantive robustness plot for all four effects: the effect of Republican and Democratic constituents on Republican legislators and on Democratic legislators. These plots are shown in Figure 5.

In Clinton's analysis, Republican and Democrat constituents' ideologies have a statistically significant relationship with their Republican legislators' voting records, but only Republican constituents' ideology is related to Democrat legislators' voting records. Substantively speaking, however, the robustness of these results varies. Republican constituents' ideology seems to have a reasonably robust relationship with both Democrat and Republican legislators' voting records (Figures 5a and 5c): even researchers who value false positives 36 times more than false negatives ( $\gamma = 6$ ) would act on the belief that the coefficient  $\lesssim 1$ . With a coefficient of 1, a one standard deviation change in constituent ideology is predicted to move a Democratic legislator at the 50th percentile of his/her party's ideology moving to the 57th percentile. The effect on Republican legislators is smaller, but still politically significant: a one standard deviation change in Republican constituents' ideology is predicted to move a Republican legislator at the 50th percentile of his/her party's ideology to the 53rd percentile. Visually, we can see that these results are substantively robust because the range of acceptable coefficients is large across a very wide swath of loss aversion coefficients  $\gamma$ .

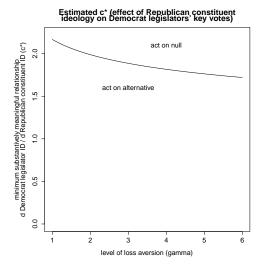
The relationship between Democrat constituents' ideology and legislators' voting records is much less robust. This may not be immediately clear from looking at the coefficients in the central purpose of the present article.

Figure 5: Substantive robustness assessment for constituent influence on "key votes" in Congress, from Clinton (2006)

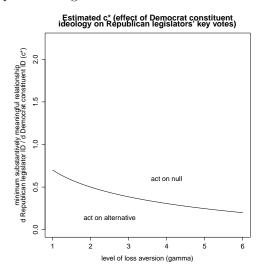
(a) Republican Constituents' Influence on Republican Legislators



(c) Republican Constituents' Influence on Democrat Legislators



(b) Democratic Constituents' Influence on Republican Legislators



(d) Democratic Constituents' Influence on Democrat Legislators

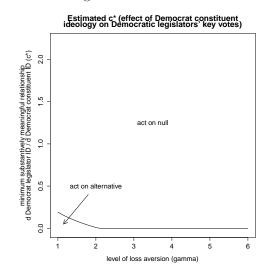


Table 1. For instance, the 90% confidence interval for Democratic constituents' effect on Republican legislators is [0.1, 1.3]. This estimate indicates considerable uncertainty about this relationship, but it also firmly excludes the null and includes the possibility of a subtantively large relationship. But as Figure 5b shows, a community of heterogeneous researchers would tend to agree on a coefficient of less than or equal to about 0.2—quite close to the left boundary of the 90% confidence interval. A coefficient of this size implies that one standard deviation change in Democrat constituent ideology moves a Republican legislator at the 50th percentile of his/her party's ideology to the 51.8th percentile. That is a reasonably small relationship, probably politically ignorable. We conclude that the relationship between Democrat constituents and Republican legislators is not robust to different standards of substantive judgment. Visually, the non-robustness of these relationships is visible because the range of acceptable coefficient sizes is close to zero (or actually zero) for even modest values of loss aversion  $\gamma$ .

Given our earlier simulation results, indicating that false positives are minimized by accepting only those results that are substantively robust and statistically significant, we are led to conclude that Clinton's evidence does not support a link between Democrat constituents' ideology and the voting records for legislators of either party. Indeed, when estimating an errors-in-variables regression on the same data, Clinton comes to the same conclusion (but on the basis of statistical significance).

### Applied example: Brambor, Clark and Golder (2006)

In an influential paper, Brambor, Clark and Golder (2006, hereafter BCG) describe the importance of interpreting interaction terms substantively, and provide tools and techniques for informative interpretation. These tools and techniques are readily adapted to substantive robustness assessment using our technique. In one example, BCG present a graphical method of presenting marginal effects from linear models with interaction terms. Their example comes from Clark and Golder (2006), a model of the relationship between presidential

Table 2: The impact of presidential elections on the effective number of electoral parties (Replication of Table 1 from Brambor, Clark and Golder, 2006)

Regressor	β	s.e.
Election Proximity	-3.53	0.54
Presidential Candidates	0.33	0.17
Proximity * Pres. Cand.	0.84	0.23
Controls	_	_
Constant	3.11	0.33

OLS model of the number of electoral parties in the legislature. Standard errors are clustered on country.  $R^2 = 0.25$ , N = 602.

elections and legislative fragmentation (the number of parties in the legislature). The model is listed in Table 2; note that the relationship between the proximity of presidential elections and the number of electoral parties in the legislature is contingent on the number of presidential candidates.

To test hypotheses about the marginal effect of presidential elections on legislative fragmentation, BCG recommend constructing a plot like the one in Figure 6a, which shows this marginal effect and its 95% confidence interval for different numbers of presidental candidates. Mathematically, the marginal effect shown in the figure is:

$$\frac{\partial \# \text{ candidates}}{\partial \text{ presidential elections}} = \beta_{proximity} + (\beta_{prox*cand} * \# \text{ candidates})$$

and therefore a combination of multiple coefficients from the regression. In Figure 6b, we augment this assessment of statistical significance with a substantive robustness assessment. Specifically, we use our software to determine  $c^*$  for each number of presidential candidates using the loss averse utility function from equation 5 and a variety of different loss aversion coefficients  $\gamma$ .<sup>13</sup> The shaded regions indicate the marginal effect sizes that could be acted upon by researchers with the corresponding level of loss aversion  $\gamma$ . The wider the shaded range, the more robust is the judgment of substantive significance.

<sup>&</sup>lt;sup>13</sup>Inference about the substantive significance of a marginal effect  $(dy/dx|z) = (y'_x|z)$  in an interaction

As the figure shows, when the number of presidential candidates is small (between 0 and  $\approx 1.5$ ) there is a very robust negative relationship between presidential elections and the number of electoral parties. Even researchers who are extremely averse to false positives ( $\gamma > 20$ , valuing a false positive 400 times more than a false negative<sup>14</sup>) should conclude that the number of legislative parties participating in an election shrinks by  $\approx 1$  when concurrent presidential elections are held (compared to legislative elections held precisely in between presidential terms, such as a U.S. midterm election). To reduce the interpretation to simple visual terms, there is a broad swath of the marginal effect space that is shaded at a very high level of  $\gamma$ , indicating that even researchers extremely averse to false positives would act on the alternative hypothesis of a large marginal effect. In short, researchers with a wide variety of standards for the interpretation of evidence would accept a large relationship between presidential elections and the number of electoral parties in the legislature.

This negative relationship becomes less robust as the number of presidential candidates grows. When there are three presidential candidates, for example, researchers who are even moderately averse to false positives ( $\gamma > 4$ ) would act on the null that there is no relationship at all. Those who are reasonably indifferent between Type I and II errors ( $\gamma \leq 2$ , valuing false positives 4 times as much as false negatives) would be willing to act on the belief that a reduction of  $\approx 1$  party occurs when presidential elections are concurrent. These same researchers would also accept a that positive relationship exists between the number of legislative parties and the fact of concurrent presidential elections when there are a very large number of presidential candidates competing for office ( $\approx 6$ ). In visual terms, only a term context makes use of samples from the distribution  $f(y'_x|f(\beta|\text{data}),z)$  to calculate the root in c of:

$$\int \left[ u(\operatorname{accept}|y_x'(\beta, z), c) - u(\operatorname{reject}|y_x'(\beta, z), c) \right] *$$

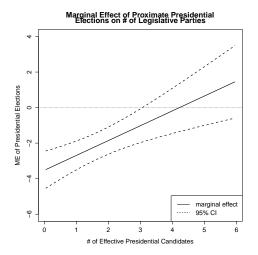
$$f(y_z'|f(\beta|\operatorname{data}), z) dy_z' = 0$$
(8)

With a sufficient number of samples from  $f(y_z'|f(\beta|\text{data}),z)$  to allow for accurate kernel density estimation of the underlying distribution, the integral on the right hand side of equation 8 can be numerically calculated as easily as that of equation 3 using standard numerical integration packages. Our software allows the calculation of a  $c^*$  value when fed a large number of samples from  $f(y_z'|f(\beta|\text{data}),z)$  and a  $\gamma$  value.

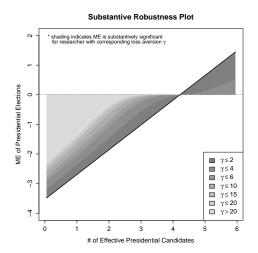
<sup>&</sup>lt;sup>14</sup>To put the preference into perspective, a person with these preferences would refuse to pay \$2.50 for a lottery ticket with a 50% chance of paying off \$1000!

Figure 6: Robustness Assessment for the Marginal Effect of Proximate Presidential Elections on Number of Electoral Parties in the Legislature (Figure 3 from Brambor, Clark and Golder, 2006)





(b) Substantive Robustness Assessment



small portion of the marginal effect space is shaded, and only for comparatively small  $\gamma$ ; this indicates that many researchers would *not* conclude that there is a substantively significant relationship (i.e., they would act on the null).

The pattern of substantive findings—a negative relationship between legislative parties and proximate presidential elections when the number of presidential parties is low, and no relationship otherwise—is consistent with the statistical significance findings depicted in Figure 6a, and with Clark and Golder's theoretical hypothesis. Based on our earlier simulations, we therefore have a great deal of confidence in the robustness of this finding.

### Conclusion

Our paper has focused on making three claims:

1. Substantive significance is easy to subjectively explain (a result strong and certain enough to justify acting upon), but hard to objectively specify because of reasonable differences in scientific judgment.

- 2. As a consequence of (1), it is more productive to ask "would researchers with widely varying preferences consider this result substantively significant?" rather than "is this result substantively significant?" The first question is equivalent to asking, "is this result substantively robust?"
- 3. Requiring that results be substantively robust as well as statistically significant results in more powerful hypothesis tests and fewer erroneous results.

It remains for us to argue for why should researchers use our particular technique to assess the substantive robustness of their results. It is both possible and useful to think about the substantive significance of results in more heuristic ways. For instance, Achen (1982) recommends examining the bounds of a 95% confidence interval of a marginal effect, thinking through the political significance of the relationship if the effect were on either boundary. We certainly do not want to argue against the value of this kind of analysis.

We think that our procedure is a valuable supplement to heuristics like those proposed by Achen (1982) because it clarifies the role of heterogeneous researcher standards in the formation of substantive judgments. By applying our technique, it becomes easier to distinguish results that are reasonably robust in one way or another—either most people would agree that a relationship exists, or that it doesn't, based on the available evidence—from those that are more contentious and for which more information is required. From a hypothesis testing perspective, the formality of our procedure also allows us to establish an important benefit of substantive robustness assessment: requiring a relationship to be statistically significant and (minimally) substantively robust improves our ability to distinguish true positives from false positives.

Going a bit further, we also think that our procedure is more rational and consistent than heuristic assessment because it is explictly built on settled, existing work in statistical decision theory and rational choice theory. It ensures that results are objectively evaluated inside of the reasonable bounds of subjective differences in judgment. And, we hope that it is no less convenient than heuristic substantive assessment because the process is preprogrammed for use in a wide variety of models and can be performed with only a few keystrokes.

Our argument does not hinge on the choice of any particular utility function to be used inside of our framework, and indeed we cannot hope to evaluate the infinite space of possible alternatives. We have shown that our particular choice is able to accommodate widely varying opinions about (a) the proper tradeoff between the risks of false positive and false negative and (b) the minimum coefficient size needed for substantive significance. We have also shown that a statistically significant result that is substantively robust under our particular utility function is extremely unlikely to be a false positive, and that requiring substantive robustness alongside statistical significances improves power (true positive detection) at any fixed size (false positive rate). We believe that is prima facie evidence for the value of our choice of utility function; it improves on the status quo of statistical significance testing only (or significance testing with informal arguments about substantive significance). But we must accede that there could be other possibilities that would be even better on these dimensions, and leave this exploration to future work. Our main point is that inferences should be robust to disagreement about important dimensions of judgment, not sensitive to them.

We also believe future work should explore certain applications that would be especially helpful to applied researchers. In particular, not all model products are simple coefficient distributions or even affine functions of coefficients (as in our illustrations). There exist, for example, joint F-tests to simulaneously test for the joint statistical significance of a group of coefficients. It would undoubtedly be helpful to assess the joint substantive significance of a group of coefficients as well. This would allow a researcher to ask whether a group of marginal effects is collectively influential enough on a dependent variable to be substantively actionable even if any individual effect is not.

#### References

- Achen, Chris. 1982. Interpreting and Using Regression. Quantitative Applications in the Social Sciences Sage University Press.
- Brambor, Thomas, William Clark and Matt Golder. 2006. "Understanding Interaction Models: Improving Empirical Analyses." *Political Analysis* 14:63–82.
- Clark, WR and Matthew Golder. 2006. "Rehabilitating Duverger's theory." Comparative Political Studies 39(6):679–708.
  - URL: http://cps.sagepub.com/content/39/6/679.short
- Clinton, Joshua. 2006. "Representation in Congress: Constituents and Roll Calls in the 106th House." *Journal of Politics* 68(2):397–409.
- Clinton, Joshua, Simon Jackman and Douglas Rivers. 2004. "The Statistical Analysis of Roll Call Data." American Political Science Review 98:355–370.
- DeGroot, Morris. 2004 {1970}. Optimal Statistical Decisions. Wiley Interscience.
- French, Simon and Diavid Rios Insua. 2000. Statistical Decision Theory. Kendall's Library of Statistics Oxford University Press.
- Gelfand, Alan E. and Malay Ghosh. 2000. Generalized Linear Models: A Bayesian View. In *Generalized Linear Models: A Bayesian Perspective*, ed. Sujit K. Ghosh and Bani K. Mallick. Marcel Dekker chapter 1, pp. 3–22.
- Gelman, Andrew, John B. Carlin, Hal S. Stern and Donald B. Rubin. 2003. *Bayesian Data Analysis, Second Edition*. 2 ed. Chapman & Hall/CRC.
- Gill, Jeff. 1999. "The Insignificance of Null Hypothesis Significance Testing." *Political Research Quarterly* 52(3):647–674.

Hunter, John E. 1997. "Needed: A Ban on the Significance Test." *Psychological Science* 8:3–7.

Lunt, Peter. 2004. "The Significance of the Significance Test Controversy: Comments on 'Size Matters'." The Journal of Socio-Economics 33:559–564.

Manski, Charles F. 2007. *Identification for Prediction and Decision*. Harvard University Press.

McCloskey, Deirdre. 1998. The Rhetoric of Economics. University of Wisconsin Press.

Miller, Jane E. 2008. Interpreting the substantive significance of multivariable regression coefficients. In 2008 Proceedings of the American Statistical Association, Statistical Education Section. URL: http://policy.rutgers.edu/faculty/miller/2008regression\_coefficients.pdf.

Pratt, John W., Howard Raiffa and Robert Schlaifer. 1996. *Introduction to Statistical Decision Theory*. MIT Press.

Siegfried, Tom. 2010. "Odds Are, It's Wrong: Science fails to face the shortcomings of statistics." Science News 177:26.

Wald, Abraham. 1950. Statistical Decision Functions. Wiley.

Ziliak, Steven and Deirdre McCloskey. 2004. "Size Matters: The Standard Error of Regressions in the American Economic Review." The Journal of Socio-Economics 33:527–546.

Ziliak, Steven and Deirdre McCloskey. 2008. The Cult of Statistical Significance: How the Standard Error Costs Us Jobs, Justice, and Lives. University of Michigan Press.