

Statistical Analysis of Endorsement Experiments: Measuring Support for Militant Groups in Pakistan

Will Bullock, Kosuke Imai, and Jacob N. Shapiro

Department of Politics, Princeton University, Princeton, NJ 08544

e-mail: kimai@princeton.edu (corresponding author)

Political scientists have long been interested in citizens' support level for such actors as ethnic minorities, militant groups, and authoritarian regimes. Attempts to use direct questioning in surveys, however, have largely yielded unreliable measures of these attitudes as they are contaminated by social desirability bias and high nonresponse rates. In this paper, we develop a statistical methodology to analyze endorsement experiments, which recently have been proposed as a possible solution to this measurement problem. The commonly used statistical methods are problematic because they cannot properly combine responses across multiple policy questions, the design feature of a typical endorsement experiment. We overcome this limitation by using item response theory to estimate support levels on the same scale as the ideal points of respondents. We also show how to extend our model to incorporate a hierarchical structure of data in order to uncover spatial variation of support while recouping the loss of statistical efficiency due to indirect questioning. We illustrate the proposed methodology by applying it to measure political support for Islamist militant groups in Pakistan. Simulation studies suggest that the proposed Bayesian model yields estimates with reasonable levels of bias and statistical power. Finally, we offer several practical suggestions for improving the design and analysis of endorsement experiments.

1 Introduction

Political scientists have long been interested in measuring sensitive political attitudes. Yet, asking sensitive questions directly in surveys has often resulted in bias due to social desirability and high item nonresponse. To overcome these problems, researchers have relied upon survey methodologies that are specifically designed for eliciting truthful responses through the means of indirect questioning. Two most commonly used techniques are randomized response methods and list experiments (see e.g., Blair and Imai 2010; Gingerich 2010; Imai 2011).

In this paper, we focus on an alternative survey methodology called *endorsement experiments*, which have been used to measure the levels of political support for socially sensitive actors such as ethnic minorities, militant groups, and authoritarian regimes.¹ In endorsement experiments, randomly selected respondents are asked to express their opinion about several policies endorsed by a socially sensitive actor of interest. These responses are then contrasted with those from a control group that receives no endorsement. If the endorsement by a political actor induces more support for policies, then this is taken as evidence for the existence of support for that actor. The main advantage of endorsement experiments is that indirect questioning may increase truthful responses, improve response rates, and enhance safety for enumerators and respondents in the case of extremely sensitive topics. The drawback, on the other hand, is that it can only provide an indirect measure of support.

Authors' note: All of the code necessary to reproduce the results in this paper can be downloaded from the Dataverse as Bullock, Imai, and Shapiro (2011). A previous version of this paper was circulated as "Measuring Political Support and Issue Ownership Using Endorsement Experiments, with Application to the Militant Groups in Pakistan."

¹See Cohen (2003), Kam (2005), Fair, Malhotra, and Shapiro (2009), and Nicholson (2011) for recent applications in political science. Sniderman and Piazza (1993) also use a similar experiment with endorsements by Congress and community organizations.

Despite the increasing popularity and promise of endorsement experiments, the common statistical methods used to analyze endorsement experiments (e.g., difference in means and regression models) have been unsatisfactory. In particular, the standard methods fail to properly combine responses across multiple policy questions, the design feature of a typical endorsement experiment. To obtain robust estimates of support levels and improve statistical efficiency, most researchers use multiple policy questions. Yet, there exists no principled way of combining responses to these different questions to produce a single measure of support level. Typical endorsement experiments also attempt to measure support levels for multiple groups. This poses another statistical challenge because one must deal with a relatively small sample size per group. Together with the loss of information due to indirect questioning, this calls for a statistical method that allows researchers to efficiently estimate levels of support for multiple groups by combining responses to several policy questions.

To overcome these methodological challenges, we develop a Bayesian measurement model to analyze such endorsement experiments.² Our model is based on item response theory, which has been used in political science to estimate ideal points of legislators from roll call data (e.g., Poole and Rosenthal 1985; Heckman and Snyder 1997; Clinton, Jackman, and Rivers 2004).³ The proposed model provides estimates of political support measured on the same scale as the ideal points of respondents by combining responses to multiple policy questions. This framework facilitates a meaningful interpretation of the magnitude of support measures.⁴

Furthermore, our model can incorporate individual level covariates to explore the association between respondents' characteristics and the degree of their support for a political actor. Thus, the proposed methodology can efficiently recover the latent support level while directly incorporating all the estimation uncertainty.⁵ We also show that our model can accommodate a multilevel structure of data via hierarchical modeling and discover spatial variation of support levels (Gelman and Hill 2007). This means that the proposed model can incorporate aggregate-level (e.g., geographical units) covariates as well as individual characteristics. Finally, we demonstrate how to conduct poststratification to further obtain efficient aggregate-level estimates of support based on census data whenever such information is available (Park, Gelman, and Bafumi 2004; Lax and Phillips 2009).

The proposed methodology clarifies the key assumptions that enable researchers to interpret the results of endorsement experiments as measures of support levels. This in turn makes it possible for us to offer practical suggestions for improving the design and analysis of endorsement experiments so that the credibility of these assumptions can be increased. Our specific recommendations for empirical researchers are summarized in the concluding section of this paper.

We offer both empirical and simulation evidence demonstrating the applicability and statistical efficiency of the proposed methodology. On the empirical side, we apply our model to a recent large-scale endorsement experiment in Pakistan and estimate levels of support among citizens for different militant groups such as al-Qa'ida. We show that the proposed method yields efficient estimates of support for each group within divisions, which are relatively small geographical units of Pakistan. The efficiency gains are especially valuable in this context as they allow us to study spatial variation in support for militant organizations across more refined political boundaries than is possible with standard methods.

Our analysis reveals interesting patterns of support for militant groups, which are not obtained by standard regression models. First, attitudes toward militant groups are quite geographically clustered. Even within the four provinces of Pakistan, there are clear clusters of support and enmity toward these groups, with Gujranwala in northern Punjab emerging as the most supportive region for most groups. This is in marked contrast to a policy discourse that focuses on Bahawalpur, in southern Punjab, as the

²All of the code necessary to reproduce the results in this paper can be downloaded from the Dataverse as Bullock, Imai, and Shapiro (2011). The replication code also can be altered for applications to other studies.

³Although some have used the item response theory to analyze surveys (e.g., Clinton and Lewis 2008; Bafumi and Herron 2010), to the best of our knowledge, none have applied it to survey experiments.

⁴The model can also yield an estimated level of support for each political actor given any particular policy rather than averaging across all policies.

⁵This contrasts with the standard two-step estimation procedure where the levels of support are first estimated and then these estimates are regressed on individual covariates of interest.

locus of support for militancy in Pakistan. Second, militant groups tend to receive less support in the areas where they conduct the most attacks. In particular, Pakistanis' support for the Afghan Taliban is almost inversely proportional to their distance from the Afghan border. This finding raises a subtle but important point that high levels of militant activity in an area should not be taken as evidence for the existence of support for militant groups. Third, we find that the individual-level variables that are the focus of much policy attention—education and income, for example, are frequently cited as key correlates of support for Islamist militancy—matter little once regional differences are taken into account. It is thus local political attitudes that are the key determinate of support, a finding with strong implications for policies to combat the threat of militancy in Pakistan.

On the simulation side, we show that our Bayesian measurement model has good (frequentist) statistical properties. Our first set of simulation studies are based on the data generating process that closely follows the Pakistan data which have a relatively large sample size. In this setup, we show that the proposed model can recover the true parameter values with little bias and strong statistical power. In addition, our Bayesian confidence intervals have a coverage rate close to the nominal rate. The second set of simulation studies uses a similar setup but is based on sample sizes of typical surveys. Here, we find that the bias of the resulting estimates remains small and that the statistical power decreases, yet stays at a reasonable level.

The rest of the paper is organized as follows. In Section 2, we briefly describe the recent large-scale endorsement experiment conducted in Pakistan, to which we apply our proposed methodology. In Section 3, we introduce our proposed Bayesian hierarchical measurement model for endorsement experiments and formally define quantities of interest. We also discuss issues related to estimation and inference based on this model. In Section 4, we apply our method to the Pakistan data and present the empirical results. In Section 5, we show the results of simulation studies and evaluate the statistical properties of the proposed Bayesian measurement model. Finally, in Section 6, we offer concluding remarks and practical suggestions for improving the design and analysis of endorsement experiments.

2 Measuring the Support for Militant Groups in Pakistan

2.1 Background

Militant violence in Pakistan stands at the top of the international security agenda, yet little is known about who supports militant organizations and why. As a result, discussions about why Pakistan suffers so much political violence tend to turn to untested assertions that poverty, poor education, and resistance to Western values drive support for militant organizations (see Shapiro and Fair, 2010, for a summary of these arguments). It is hard to overstate the influence these ideas have. United States and Western policies toward Pakistan have devoted billions of dollars to encouraging economic and social development as an explicit means of diminishing the militant threat.

In some ways, this state of affairs is not surprising. Directly asking respondents how they feel about militant organizations is problematic, to say the least, in places suffering from political violence. In particular, there are serious safety concerns for enumerators and respondents who discuss such sensitive issues. Item nonresponse rates are often quite high as respondents understandably fear that providing the “wrong” answer will lead to social censure and may even threaten their safety. Moreover, direct questions are subject to social desirability bias insofar as answers combine respondents' true attitudes with what they believe to be the socially appropriate response.

One way to overcome these problems is through the use of an endorsement experiment. This approach measures the differences in support for various policies between two groups: one group of respondents who are told only about the policy and another group who are told that a militant organization endorses it. If respondents are randomly divided into these two groups, then the differences between the two conditions reveal how much support for a policy changes by being associated with a militant group, thereby providing an indirect measure of support for that group. Unlike a direct measure, nonresponse and social desirability biases are minimized since respondents are reacting to the policy and not directly to the group itself. By asking respondents about multiple policy issues and randomizing the pairing of issue with group, we can identify effects for multiple groups that are unlikely to be influenced by the details of any specific policy.

2.2 The Pakistan Survey Experiment

Fair, Malhotra, and Shapiro (2009) implement such a design to study support for four groups—Pakistani militants fighting over Kashmir, militants fighting in Afghanistan (a.k.a. the Afghan Taliban), al-Qa'ida, and the sectarian militias or *firqavarana tanzeems*—using a large nationally representative sample of Pakistanis. The survey entailed two main innovations. First, as noted, the researchers used an endorsement experiment to avoid the high item nonresponse rates and social desirability bias that plagued previous efforts to study the politics of militancy in Pakistan. Second, the researchers employed a sample that permitted them to make reliable inferences about a range of political attitudes in each of the four provinces of Pakistan: Punjab, Sindh, North Western Frontier Province (NWFP), and Balochistan. The name of NWFP was later changed to Khyber-Pakhtunwa province, but in this paper we use NWFP.

Working with Pakistani partners, Socio-Economic Development Consultants (SEDCO), the researchers drew a sample of 6000 adult Pakistani men and women. The respondents were selected randomly within 500 primary sampling units (PSU). Following the rural/urban breakdown in the Pakistan census, the survey entailed 332 rural PSU and 168 urban ones. Because Pakistan's provinces are heavily skewed—the NWFP and Balochistan only account for 14.5% and 4.9% of Pakistan's population, respectively (Pakistan Federal Bureau of Statistics 2008, 68)—the survey over-sampled these smaller provinces. The face-to-face questionnaire was fielded by six mixed-gender teams between April 21, 2009 and May 25, 2009.

Perhaps in part due to the indirect questioning, the overall response rate exceeded 90%, approaching the response rates achieved by the United States Census Bureau. The Pakistan endorsement experiment had the following features. First, respondents were randomly assigned to “treatment” or “control” groups (one half of the sample was assigned to each group). Respondents in the control group were asked their level of support for four policies:

- World Health Organization (WHO) plan to provide universal polio vaccinations in Pakistan.
- Government of Pakistan plan for curriculum reforms in religious schools or *madaris*.
- Reforming the Frontier Crimes Regulation (FCR) governing the tribal areas.
- Using peace jirgas to resolve disputes over the Afghan border, the Durand Line.

Respondents in the treatment group were asked identical questions but were told that one of the four groups supports the policy in question. The group associated with each of the four policies was randomized within the treatment group, effectively randomizing the order of the endorsements.

For example, the script for the control group used in this survey experiment reads as follows for the WHO polio vaccination plan question.

The World Health Organization recently announced a plan to introduce universal polio vaccination across Pakistan. How much do you support such a plan?

- (1) A great deal; (2) A lot; (3) A moderate amount; (4) A little;
- (5) Not at all.

In contrast, the script for a treatment group with respect to the same policy question reads as follows:

The World Health Organization recently announced a plan to introduce universal polio vaccination across Pakistan. Pakistani militant groups fighting in Kashmir have voiced support for this program. How much do you support such a plan?

- (1) A great deal; (2) A lot; (3) A moderate amount; (4) A little;
- (5) Not at all.

These four policies were chosen in consultation with SEDCO and other Pakistani colleagues to meet three criteria. First, each policy issue was being discussed in Pakistan at the time. Second, an informed citizen would be likely to know about these policies. Third, citizens' opinions over the policies were not rigid so that the group endorsements could have an effect. All three criteria were met based on the results of a 200-person pretest and four focus groups conducted during the survey design phase (see also Section 6 for practical suggestions about how to improve the design of endorsement experiments).

Figure 1 graphically displays the aggregate distribution of responses across four provinces from the Pakistan survey experiment. Pakistanis held generally favorable views toward each of the four policies, with greatest support for education reforms and polio vaccinations. The plots show that there exist significant differences across provinces (four rows of the graph) for each policy question. On the other hand, in each province, the aggregate distribution of responses does not vary dramatically between the control group and each treatment group for any of the four policy questions. Thus, the main methodological challenge is whether we will be able to detect any systematic pattern of support levels from these data using the proposed method.

2.3 Substantive Issues

Before turning to a detailed discussion of the proposed method, we briefly highlight key substantive issues we hope to address in our empirical analysis. First, as with all political organizations, militant groups have highly localized components. Understanding the patterns of support in more detail can help identify the specific grievances that drive people to embrace these organizations. Second, in addition to the spatial variation of support levels for militant groups, researchers are interested in how individual characteristics—gender, education, and socioeconomic status, for example—are associated with support levels after taking into account regional differences. This is not merely an academic concern. Foreign aid programs and other policies intended to combat militancy often rest on implicit hypotheses about which individuals are most susceptible to militants' appeals. (e.g., United States Agency for International Development 2009). Two of the most common such hypotheses are that (a) the poor and (b) the uneducated are more likely to support and/or participate in violent political organizations. As we will see, there is little evidence for either hypothesis in these data once regional differences are accounted for.

Finally, we also investigate the relationship between levels of political violence and support for militant groups. There is striking variation in levels of violence across Pakistan. In the year before the survey was fielded, for example, the number of people killed in terrorist attacks ranged from 11,429 in NWFP (population of 17 million) to 1791 in Sindh province (population of 30 million), almost a ten-fold difference in per capita violence. And within each province, there exists substantial further variation in violence across smaller political units. Understanding how levels of violence correlate with support for militant groups can provide vital insight into how militants interact with the population at large. If support is higher in more violent areas, it could suggest that popular approval is a key enabler for violent groups. If, however, support is lower in more violent areas, that could suggest that the negative externalities of terrorist violence may actually alienate the population, at least locally.

Answering all these substantive questions requires researchers to obtain precise estimates of support within relatively small geographical boundaries. Standard methods of analyzing endorsement experiments cannot do so without prohibitively large samples. We therefore turn to the proposed methodology, which allows us to combine responses to multiple policy questions in a principled manner and obtain efficient estimates of support levels for several militant groups at both individual and aggregate levels.

3 The Proposed Methodology

In this section, we introduce a new methodology for analyzing endorsement experiments. We begin by formally describing the design of the survey experiment and then introduce the Bayesian hierarchical measurement model of political support. Under this model, we define quantities of interest and briefly explain the issues related to estimation and inference.

3.1 The Basic Design of Endorsement Experiments

Suppose that we wish to measure the level of support for K political actors within a target population \mathcal{P} . To do so, we consider a survey experiment for a random sample of N respondents from \mathcal{P} . In the survey, respondents are asked whether or not they support each of J policies chosen by researchers. Here, for the sake of simplicity, we begin by assuming that respondents are answering a binary question regarding their support for a policy. We later generalize our model to ordinal responses so that respondents can be asked to specify the degree of support for a policy using an ordinal scale, as done in the Pakistani experiment.

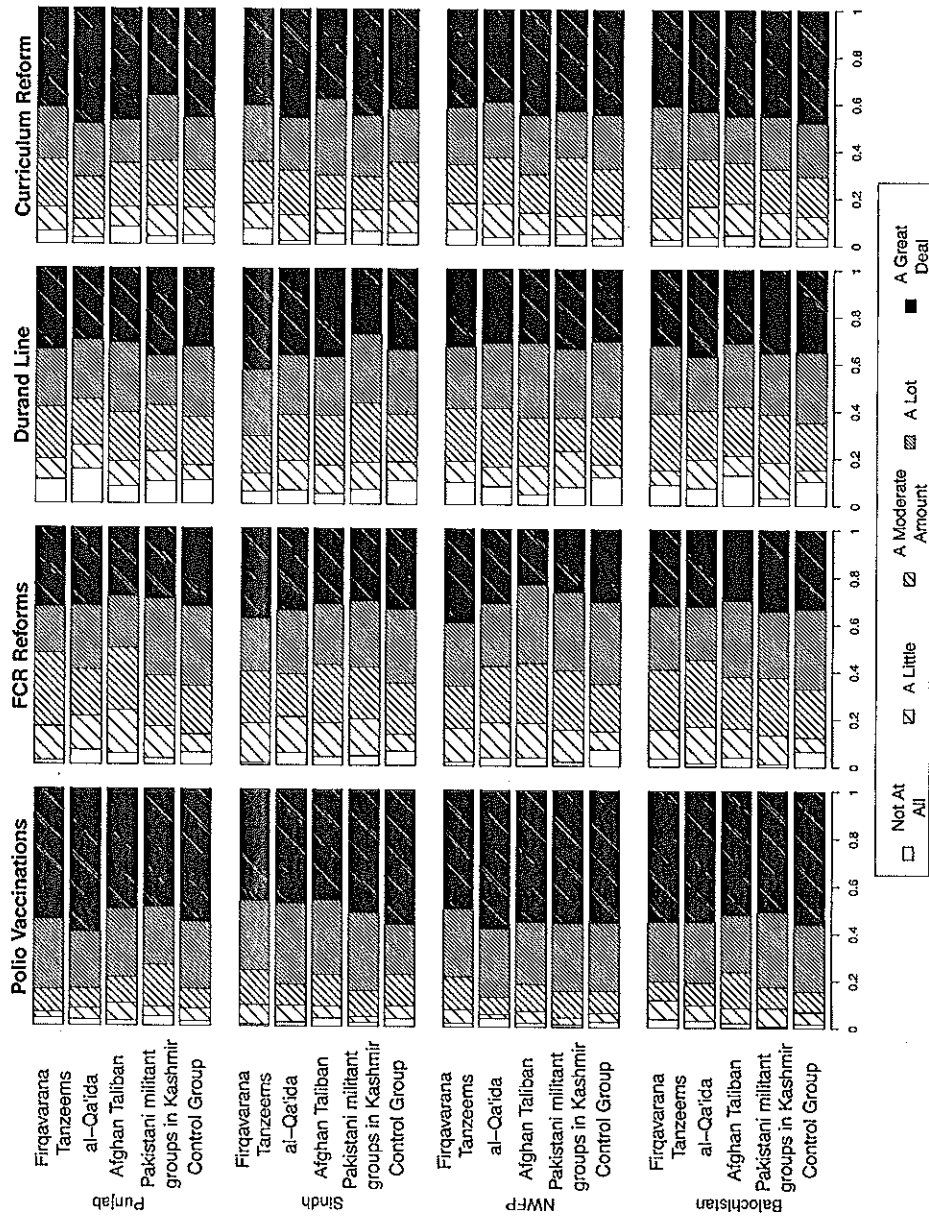


Fig. 1 Aggregate distribution of responses from the Pakistan survey experiment. The aggregate distribution of responses varies most notably across issues and provinces but does not vary dramatically between the control group and each treatment group for any of the four policy questions.

We assume that these J policies are a representative sample of a population pool of policies \mathcal{J} (i.e., a policy space), with respect to which the level of support for political actors is to be measured. If we index each policy by j , we can use Y_{ij} to represent respondent i 's answer to the survey question regarding policy j , where $Y_{ij} = 1$ ($Y_{ij} = 0$) implies that respondent i favors (opposes) policy j . We denote an M -dimensional vector of the observed characteristics of respondent i by Z_i .

The key aspect of the survey experiment is to attach the endorsement by a particular political actor to each of these policy questions. We then observe whether the level of support for these policies changes in comparison to the case where no such endorsement is given. We randomly assign one of K political actors as an endorser to respondent i 's question regarding policy j and denote this "treatment" variable by $T_{ij} \in \{0, 1, \dots, K\}$. We use $T_{ij} = 0$ to represent the "control" observations where no political endorsement is attached to the question. Alternatively, one may use the endorsement by a neutral actor as the control group.

Using the standard statistical framework of causal inference (Holland 1986), this implies that we can define $(K + 1)$ potential outcomes for respondent i 's answer to the question regarding policy j . Thus, $Y_{ij}(0)$ represents respondent i 's potential answer without any political endorsement and $Y_{ij}(k)$ is respondent i 's potential answer when policy j is endorsed by political actor k . Although we only observe one of these potential outcomes, that is, $Y_{ij} = Y_{ij}(T_{ij})$, the randomization of the treatment implies that the average treatment effect can be identified using the difference in means between treatment and control groups because of the independence between the treatment and the potential outcomes.⁶

There are a number of challenges that researchers must overcome in order to effectively use endorsement experiments. First and foremost, researchers must choose a set of policies to ask about and our analysis provides some guidance here. Other choices include question wording, question ordering, and possible carryover effects from one question to another. Although the discussion of these other design issues is beyond the scope of this paper, we emphasize that researchers must take these issues seriously before applying the proposed statistical methodology. We provide some practical suggestions about how to design endorsement experiments in Section 6.

3.2 The Bayesian Measurement Model

Given the basic design of endorsement experiments described above, we propose a Bayesian measurement model of political support. The main feature of the proposed model is to analyze survey experiments within the framework of the standard ideal point models used in political science.

3.2.1 The basic binary response model

We begin by assuming that the quadratic utility function can completely characterize each respondent's answer to the question regarding a particular policy. Following the notation of Clinton, Jackman, and Rivers (2004), let ζ_{j1} and ζ_{j0} denote the location of favoring and opposing policy j in the single-dimension policy space.⁷ Note that it is also possible to generalize this model to a multidimensional policy space, but such an extension is likely to be of little practical use given the limited number of issues typically used in endorsement experiments. As such, we recommend that researchers who want to use this framework carefully choose policy questions from a single policy dimension when designing endorsement experiments.

Given this setup, using x_i to denote respondent i 's ideal point, we can write the quadratic utility functions as,

$$U_i(\zeta_{j1}, k) = -\|(x_i + s_{ijk}^* - \zeta_{j1})\|^2 + \eta_{ij1}, \tag{1}$$

$$U_i(\zeta_{j0}, k) = -\|(x_i + s_{ijk}^* - \zeta_{j0})\|^2 + \eta_{ij0}, \tag{2}$$

where η_{ij0} and η_{ij1} represent stochastic error terms with mean zero and finite variances and s_{ijk}^* is the shift in the ideal point of respondent i induced by the endorsement of political actor k for favoring and

⁶Formally, the randomization implies $\{Y_{ij}(k)\}_{k=0}^K \perp\!\!\!\perp T_{ij}$ for all i and j , and hence, we have $\mathbb{E}(Y_{ij}(k) - Y_{ij}(0)) = \mathbb{E}(Y_{ij} | T_{ij} = k) - \mathbb{E}(Y_{ij} | T_{ij} = 0)$ for any k .

⁷See Peress and Spirling (2010) for an alternative utility model that results in the same statistical model.

opposing policy j . Therefore, s_{ijk}^* represents the degree to which political actor k can influence respondent i 's preference over policy j . Under the control condition of no political endorsement, there is no shift in the ideal point and thus, $s_{ij0}^* = 0$ for all i and j .

Under this framework, the probability that respondent i favors policy j with endorsement by political actor k is given by

$$\begin{aligned} \Pr(Y_{ij} = 1 \mid T_{ij} = k) &= \Pr(Y_{ij}(k) = 1) \\ &= \Pr(U_i(\zeta_{j1}, k) > U_i(\zeta_{j0}, k)) \\ &= \Pr(\alpha_j + \beta_j(x_i + s_{ijk}^*) > \epsilon_{ij}), \end{aligned} \quad (3)$$

for $k = 0, 1, \dots, K$, where $\alpha_j = \zeta_{j0}^2 - \zeta_{j1}^2$, $\beta_j = 2(\zeta_{j1} - \zeta_{j0})$, and $\epsilon_{ij} = \eta_{ij1} - \eta_{ij0}$. The first equality follows from the randomization of the treatment and implies that we can identify the marginal distribution of the potential responses. We assume that ϵ_{ij} is identically and independently distributed according to some distribution, such as the standard normal distribution and the logistic distribution.

The above formulation clarifies the key identification assumption necessary for inferring support levels from endorsement experiments. Specifically, endorsements are assumed to have no influence on respondents' interpretations of policy questions (relative to their interpretations in the absence of such endorsements). Formally, this assumption is incorporated into the model by removing subscript k from the question related parameters, α_j and β_j . If this assumption does not hold, the changes in responses induced by endorsements may reflect differences in question interpretations rather than support for endorsers. Although this assumption can neither be relaxed nor tested, we discuss how to design endorsement experiments in order to make the assumption more plausible (see Section 6).

Within this framework, we define respondent i 's level of support for political actor k given policy j such that a positive value of this parameter implies that the respondent favors the political actor and thereby a greater probability that the respondent prefers the policy. Formally, this support parameter is defined as,

$$s_{ijk} = \begin{cases} s_{ijk}^* & \text{if } \beta_j \geq 0 \\ -s_{ijk}^* & \text{otherwise.} \end{cases}$$

Here, the quadratic utility function is convenient because regardless of respondent i 's ideal point x_i and the error distribution, a greater value of this support parameter s_{ijk} implies a greater probability that the respondent favors policy j when endorsed by group k .⁸ Intuitively, if a viewpoint favoring policy j is located ideologically left of a viewpoint opposing the policy (i.e., $\zeta_{j1} < \zeta_{j0}$ and thus $\beta_j < 0$), the ideal point of a respondent who supports a political actor should be shifted to the left (i.e., $s_{ijk}^* < 0$ and thus $s_{ijk} = -s_{ijk}^*$) if the policy is endorsed by this actor.

In many situations, researchers are interested in modeling how respondents' levels of support for each political actor change as a function of their characteristics. Under our model, this can be accomplished, for example, by modeling support parameters s_{ijk} in the following hierarchical manner,

$$\begin{aligned} s_{ijk} &\overset{\text{indep.}}{\sim} \mathcal{N}(Z_i^\top \lambda_{jk}, \omega_{jk}^2), \\ \lambda_{jk} &\overset{\text{i.i.d.}}{\sim} \mathcal{N}(\theta_k, \Phi_k), \end{aligned}$$

where Z_i is an M -dimensional vector of respondent i 's covariates, λ_{jk} , and θ_k are M -dimensional vectors of unknown coefficients, ω_{jk}^2 is a unknown scalar variance, and Φ_k is a unknown $M \times M$ covariance matrix. Since the support parameter is zero for the control condition, that is, $s_{ij0} = 0$ for all i and j , we also have the restrictions $\lambda_{j0} = 0$ and $\omega_{j0} = 0$ for all j . Similarly, we can model ideal points x_i using the same covariates as follows,

$$x_i \overset{\text{indep.}}{\sim} \mathcal{N}(Z_i^\top \delta, \sigma_x^2),$$

where δ is an M -dimensional vector of unknown coefficients and σ_x^2 is a known prior variance.

⁸Formally, the sign of the first derivative of equation (3) with respect to s_{ijk}^* is the same as the sign of β_j .

3.2.2 Model interpretation

What is the key assumption that allows researchers to interpret change in support for a policy as support for an actor who endorses it? As discussed above and is clear from equations 1 and 2, our model formulation assumes that the endorsement of a policy by an actor may affect the location of respondents' ideal points but not the locations of favoring and opposing the policy. This assumption will be violated, for example, if respondents are uninformed about the policy and the endorsement provides them with additional information about the content of the policy. Under this alternative scenario, change in support for the policy can be thought of as the result of learning (different question interpretations) rather than that of expressing support for the endorser.

This ambiguity of interpretation is a fundamental problem of endorsement experiments and has little to do with what statistical method researchers employ for analysis. Unfortunately, the observed data cannot tell which interpretation is valid. One possible strategy is to ask each respondent to locate policies on the relevant policy dimension. If we find that the endorsement by an actor does not influence their location, then we may more credibly conclude that change in support for these policies can be attributed to support for the actor. Other possible strategies include choosing policies that are well known to respondents and focusing one's analysis on respondents who are well informed about these policies.

3.2.3 Prior specification

Finally, our Bayesian model is completed by placing the following independent diffuse prior distributions on all unknown parameters,

$$\begin{aligned} \alpha_j &\overset{\text{i.i.d.}}{\sim} \mathcal{N}(\mu_\alpha, \sigma_\alpha^2), \\ \beta_j &\overset{\text{i.i.d.}}{\sim} \mathcal{N}(\mu_\beta, \sigma_\beta^2), \\ \delta &\sim \mathcal{N}(\mu_\delta, \Sigma_\delta), \\ \theta_k &\overset{\text{i.i.d.}}{\sim} \mathcal{N}(\mu_\theta, \Sigma_\theta), \\ \omega_{jk}^2 &\overset{\text{i.i.d.}}{\sim} U[0, 10], \\ \text{diag}(\Phi_k) &\overset{\text{i.i.d.}}{\sim} U[0, 10], \end{aligned}$$

where Φ_k is assumed to be a diagonal matrix.

3.2.4 Generalization to the Ordinal Response Model

To measure the level of support for particular policies, many surveys use ordinal responses (e.g., strongly oppose, oppose, neutral, favor, and strongly favor) rather than binary responses. Here, we generalize the above measurement model to the ordinal response setup. We use this generalized model in our empirical application.

Suppose that the response variable Y_{ij} is the ordered factor variable taking one of L levels, that is, $Y_{ij} \in \{0, 1, \dots, L - 1\}$, where $L > 2$. We assume that a greater value of Y_{ij} indicates a greater level of support for policy j . Let ζ_{jl} denote the location of choosing the l th level of support in the single-dimensional policy space. As before, the key identification assumption is that these question parameters do not depend on endorsements. This assumption eliminates the possibility that different question interpretations (rather than support for endosers) change responses. Finally, we assume that the elements of $\{\zeta_{jl}\}_{l=0}^{L-1}$ are equally spaced in the policy space in either a descending or ascending order, that is, $\zeta_{j,l+1} - \zeta_{jl} = \psi_j$ does not vary across levels l although it may vary across policies j . This assumption can be relaxed but we do not pursue such a strategy here to keep the model parsimonious (see below for more discussion on this point).

Then, the utility function can be written as,

$$U_i(\zeta_{jl}, k) = -\|(x_i + s_{ijk}^*) - \zeta_{jl}\|^2 + v_{ijl},$$

which in turn implies the following probability model,

$$\begin{aligned} \Pr(Y_{ij} \leq l \mid T_{ij} = k) &= \Pr(Y_{ij}(k) \leq l) \\ &= \Pr(U_i(\zeta_{j,l+1}, k) < U_i(\zeta_{j,l}, k)) \\ &= \Pr(\alpha_{jl} - \beta_j(x_i + s_{ijk}^*) > \epsilon_{ijl}), \end{aligned} \quad (4)$$

for $l = 0, 1, \dots, L - 2$, where $\alpha_{jl} = \zeta_{j,l+1}^2 - \zeta_{j,l}^2 = \psi_j(\zeta_{j,l} + \zeta_{j,l+1})$, $\beta_j = 2(\zeta_{j,l+1} - \zeta_{j,l}) = 2\psi_j$, and $\epsilon_{ijl} = v_{ijl} - v_{ij,l+1}$ is assumed to be independently and identically distributed. Note that the relationship between α_{jl} and β_j implies $\alpha_{jl} \leq \alpha_{j,l+1}$ for any j and l . The resulting model then is equivalent to the *generalized partial credit model* in the psychometrics literature (see Muraki 1992). If the assumption of equal spacing between $\zeta_{j,l}$ and $\zeta_{j,l+1}$ is relaxed, then β_j will need to vary across levels l as well as across policies j (note that both the intercept α_{jl} and the slope β_j are allowed to vary across policies). However, we will not pursue this approach because there is little information in the data to estimate a separate slope parameter for each response category.

3.3 Quantities of Interest

Under the above model, our two basic quantities of interest are defined as,

$$\tau_{jk}(Z_i) = Z_i^\top \lambda_{jk}, \quad (5)$$

$$\kappa_k(Z_i) = Z_i^\top \theta_k. \quad (6)$$

Whereas $\tau_{jk}(Z_i)$ represents the average support level for political actor k with respect to policy j among individuals with characteristics Z_i , $\kappa_k(Z_i)$ denotes the average support level for the political actor among the same set of individuals without reference to any specific policy. While it is certainly of interest to examine how these two measures of support vary as a function of individual characteristics Z_i , the overall average support levels for each political actor in the population are also important, that is, $\bar{\tau}_{jk} = \mathbb{E}(\tau_{jk}(Z_i))$ and $\bar{\kappa}_k = \mathbb{E}(\kappa_k(Z_i))$, where the expectations are taken over the population distribution of Z_i .

Although the above measures facilitate the comparison of support levels across different political actors and policies, one important limitation is that their magnitude is difficult to interpret because they are measured in the latent policy space. To address this issue, we standardize levels of support in terms of the posterior standard deviation of ideal points. That is, support for political actors is measured relative to the variation in ideal points among the population. An endorsement impact of 2.0, for example, can now be interpreted as shifting a respondent's ideal point two SDs. In the U.S. context, this would be interpreted as having a moderate Republican provide responses akin to those of a moderate Democrat, or vice versa. Thus, the standardized measures provide an important intuition for interpreting their magnitude by placing them on the same scale as ideal points of survey respondents.

3.4 Identification, Estimation, and Inference

Since we have the standard ideal points model for the control group, the local identification results of Rivers (2003) are directly applicable. Following the convention in the literature, we typically constrain the distribution of estimated ideal points to have mean zero and unit variance. In addition, for global identification, we constrain the sign of at least one element of β_j for a particular policy j (see also Bafuni et al. 2005). Once the item parameters are identified from the control group, the identification of the average support parameters λ_{jk} is immediate because of randomization. The simulation studies given in Section 5 also show that all relevant parameters are estimated with little bias even in a relatively small sample.

Following many of the recent applications of the item response theory in political science, we use the Markov chain Monte Carlo (MCMC) algorithm to fit the proposed Bayesian measurement model. One can construct the MCMC algorithm using the standard data augmentation scheme with the latent outcome variable Y_{ij}^* underlying the binary or ordered response variable Y_{ij} . One important advantage of Bayesian simulation is that the estimation uncertainty can be easily calibrated using the draws from the posterior distribution.

For actual model fitting, we use the open-source software JAGS (Plummer 2009, version 2.2.0). As demonstrated in Section 5, the MCMC algorithm implemented through JAGS recovers the poster distribution well in the types of data sets we have analyzed. Our replication code is made available as Bullock, Imai, and Shapiro (2011) so that other researchers can apply the proposed methods to their own endorsement experiments.

3.5 Weighting and Poststratification

Endorsement experiments require a large sample size relative to other approaches especially when the number of political actors and/or the number of policy questions is large. To recoup this loss of efficiency, the use of effective stratification or other sampling techniques at the design stage may be important. Incorporating the resulting survey weights into estimates of support measures obtained from our model is relatively straightforward. Since s_{ijk} represents the level of support at the individual level, survey weights can be directly applied to its estimate when computing the estimated average level of support.

Poststratification at the analysis stage provides another way of estimating levels of support in the population (Park, Gelman, and Bafumi 2004; Lax and Phillips 2009). When the population joint distribution of covariates Z_i is available, this can be easily done in our model by properly weighting the estimated average levels of support, $\tau_{jk}(Z_i)$ or $\kappa_k(Z_i)$, according to the distribution of Z_i . We illustrate the use of poststratification in our analysis of the Pakistan data in Section 4.2.

4 Empirical Results

In this section, we apply the proposed methodology to the data from the Pakistan survey experiment discussed in Section 2. We consider two multilevel/hierarchical models that extend the base model described in Section 3. The first model estimates the levels of support for each militant group at the division level, where the 26 divisions of Pakistan are a now-defunct political unit that is smaller than the four provinces but sufficiently large that we have a reasonable sample in 21 of 26 divisions. The second model adds individual-level covariates to the first model to examine the variation across individuals. The two models we consider here share the same top-level specification, which is given in equation (4) where, based on substantive grounds, we assume that $\beta_j \geq 0$ for all policy questions. This assumption is easily justified for the curriculum reform question and polio vaccination question. Pakistan experts we talked to, our focus group participants, and our enumerators all expected those who support Islamist militancy to oppose them.⁹ Both models are also estimated on the set of completed surveys, $n = 5212$. For model fitting, we run three parallel chains with overdispersed starting values and monitor the convergence using the standard diagnostic statistic (Gelman and Rubin 1992). After the convergence has been achieved according to this criteria, we keep every 30th of the last 60,000 posterior draws for our inference.

4.1 Estimated Support at the Division Level

Our first model estimates the levels of support for each militant group while assuming that the degree of support is identical across policies. This leads to the following specification:

$$\begin{aligned}
 x_i &\overset{\text{indep.}}{\sim} \mathcal{N}(\delta_{\text{division}[i]}, 1), \\
 s_{ijk} &\overset{\text{indep.}}{\sim} \mathcal{N}(\lambda_{k,\text{division}[i]}, \omega_k^2), \\
 \delta_{\text{division}[i]} &\overset{\text{indep.}}{\sim} \mathcal{N}(\mu_{\text{province}[i]}, \sigma_{\text{province}[i]}^2), \\
 \lambda_{k,\text{division}[i]} &\overset{\text{indep.}}{\sim} \mathcal{N}(\theta_{k,\text{province}[i]}, \Phi_{k,\text{province}[i]}),
 \end{aligned}$$

⁹While the majority of our interlocutors expected Islamists to oppose FCR reform and jirgas over the Durand line, opinion was more diverse and some argued Islamists would favor those policies. We have therefore also run the analysis assuming $\beta_j \geq 0$ only for the curriculum reform and polio questions. The results do not change substantively for alternative assumptions.

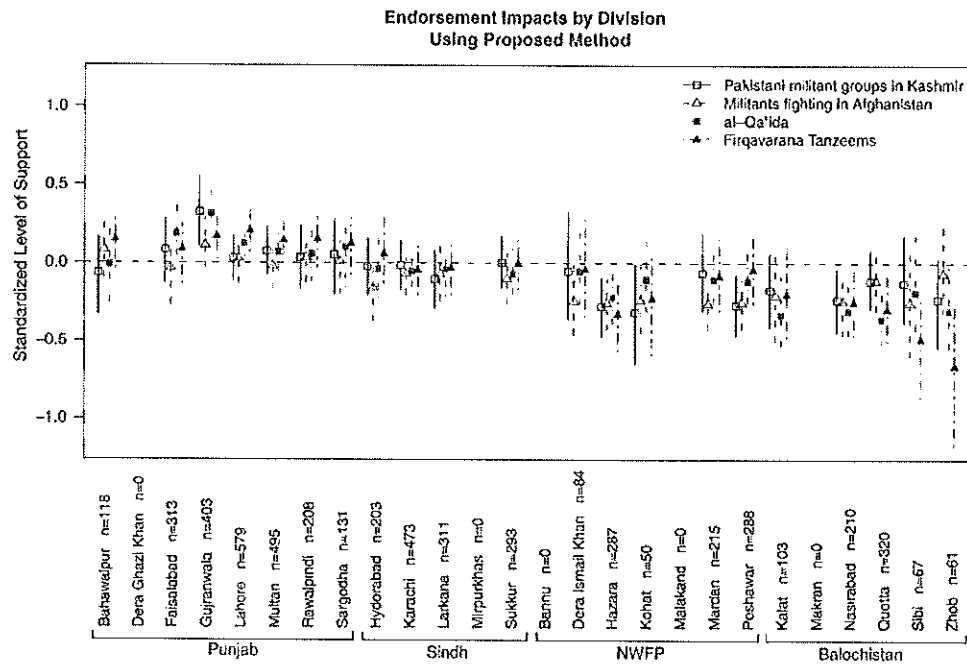


Fig. 2 Estimated division level support for each militant group. The vertical axis represents the estimated levels of support, which are standardized in terms of the posterior standard deviation of ideal points distribution. The symbols represent the point estimates and the vertical lines are 90% Bayesian confidence intervals. In this model, the degree of support is assumed to be constant across policies. The plot suggests that the levels of support varies substantially across divisions and provinces. In contrast, little variation is found across different militant groups within each division.

where partial pooling occurs across divisions within each province as well as across individuals within each division.¹⁰

Figure 2 presents the estimated levels of support for each militant group at the division level that are based on our first model. As discussed in Section 3.3, we standardize our measures of average support for each group by dividing it with the posterior standard deviation of ideal points distribution so that they are easier to interpret than the raw estimates. The plot suggests that the levels of support vary substantially across divisions within each province as well as across provinces. For example, whereas there is clearly positive support for three of four militant groups in Gujranwala (Punjab), levels of support in Hazara (NWFP) and Nasirabad (Balochistan) are estimated to be negative for all four groups. Moreover, whereas the estimated support tends to be negative within Balochistan and the NWFP, the citizens of Punjab appear to be more favorably disposed toward militant groups.

In contrast to this regional variation, however, the difference in support levels across militant groups within each division is relatively small. Although some variations exist in such divisions as Sibi and Zhob in the Balochistan province where *firqavarana tanzeems* attracts relatively large negative support, the group's 90% confidence interval overlaps substantially with those of the other militant groups.

4.2 Estimated Effects of Individual-Level Covariates and Use of Poststratification

The second model we consider is the same as the first model presented in Section 4.1 except that we incorporate individual-level covariates in order to answer substantive questions and gain statistical efficiency. Here, we use education, income, gender, and whether respondents are from an urban area.¹¹ The

¹⁰We also fitted another model that allows the intercepts α_{jt} in equation (4) to vary across divisions and place an independent and diffuse prior. The results are similar to the ones presented here.

¹¹Following the work of Fair, Malhotra, and Shapiro (2009), we code respondents' income as 1 if their household income was within the bottom quartile for their province, as 2 if their household income was within the middle quartiles for their province, and as 3

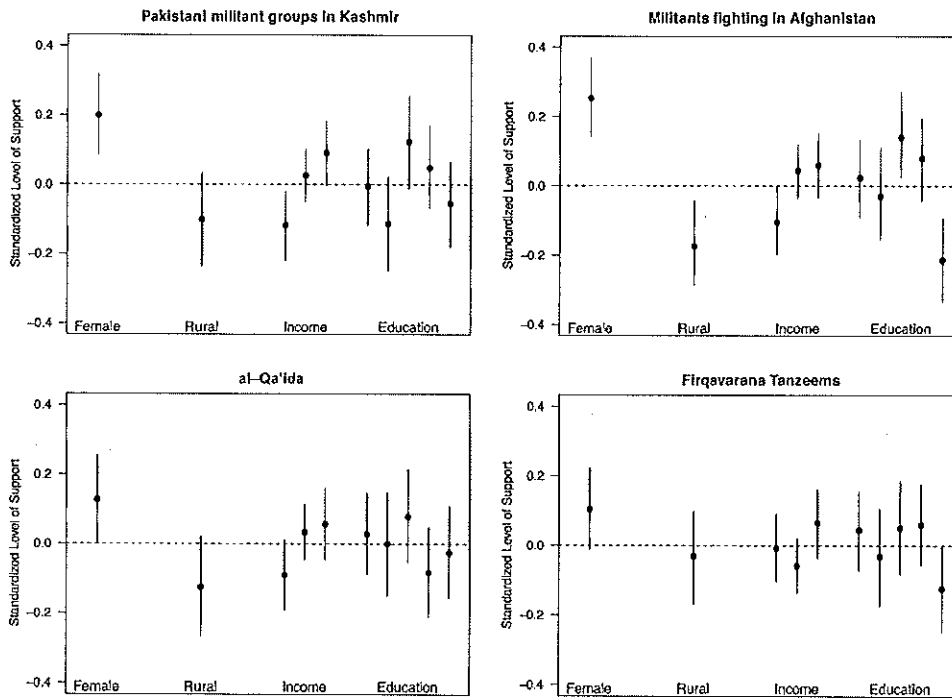


Fig. 3 Estimated effects of individual covariates on levels of support. The vertical axis represents the estimated effects of individual level covariates on levels of support, which are standardized in terms of the posterior standard deviation of ideal points distribution. The symbols represent the point estimates and the vertical lines are 90% Bayesian confidence intervals. The plot shows that after accounting for regional differences in support levels, women are slightly more supportive of each group and respondents from rural areas as marginally less supportive of each group. There appears to be no clear pattern of support across income or education groups after accounting for geographic variation.

first two variables are commonly thought to predict support for extremism, with the poor and uneducated expected to be more supportive. The latter two variables map onto a broad range of differences within Pakistani society (media consumption, employment rates, number of children, social conservatism, etc.) that are not captured in large-scale surveys. Accounting for them therefore may enhance the efficiency of our inference. In addition, since the joint population distribution of these variables is available from the recent Pakistan Social and Living Standards Measurement Survey (PSLM), we can poststratify the resulting estimates on these covariates to compute the division-level estimates of support level that adjust for remaining discrepancies between the sample and the population. This model with individual covariates are given by¹²

$$\begin{aligned}
 x_i &\overset{\text{indep.}}{\sim} \mathcal{N}(\delta_{\text{division}[i]} + Z_i^T \delta^Z, 1), \\
 s_{ijk} &\overset{\text{indep.}}{\sim} \mathcal{N}(\lambda_{k,\text{division}[i]} + Z_i^T \lambda_k^Z, \omega_k^2), \\
 \delta_{\text{division}[i]} &\overset{\text{indep.}}{\sim} \mathcal{N}(\mu_{\text{province}[i]}, \sigma_{\text{province}[i]}^2), \\
 \lambda_{k,\text{division}[i]} &\overset{\text{indep.}}{\sim} \mathcal{N}(\theta_{k,\text{province}[i]}, \Phi_{k,\text{province}[i]}).
 \end{aligned}$$

Figure 3 presents the results based on this model. The figure shows that individual covariates explain only a small portion of the variance in support levels once we take into account regional differences in

if their household income was within the top quartile for their province. We code respondents' education as 1 if illiterate, as 2 if they have only primary school education, as 3 if they have middle school education, as 4 if they have matric education, and as 5 if they have intermediate or higher education.

¹²We also fitted a model that allows the intercepts a_{jl} to vary across divisions. The results are similar to the ones we obtain here.

support across divisions. It appears that women are slightly more supportive of each group, and respondents from rural areas are marginally less supportive of each group. Contrary to common expectations, there is no monotonic relationship between education and support across groups. In addition, low-income respondents appear to dislike the three externally focused groups more than their high-income countrymen, with the difference being largest for militants fighting in Kashmir. Critically, the poor are no more or less supportive of the domestically focused *firqavarana tanzeem* than others. One interpretation of these differences is that the poor generally dislike militants as they are most exposed to the negative externalities of violence but that this tendency is tempered with respect to the *firqavarana tanzeem* who are widely understood to be executing a class-based agenda through the language of apostasy in regions where repressive land owners have historically been predominantly Shi'a.¹³

Interesting as they are, these results should be interpreted with caution given that the predictive power of these covariates is small relative to the recovered estimates from geographic divisions. As a result, the division-level estimates produced by employing poststratification (data not shown) are little changed from the estimates shown in Fig. 2. The implication is that regional politics appears to be much more important for understanding support for militancy in this context than are individual characteristics.

4.3 Substantive Discussion of Empirical Results

One way to understand the result of our Bayesian measurement model is to view the division-level support estimates overlaid on a map of Pakistan, as shown in Fig. 4. Viewing the data in this manner highlights the spatial patterns of support for militant organization that are not necessarily clear from tables and figures. Examining the spatial variation is especially important given our finding that individual covariates matter little for predicting the levels of support once regional differences are taken into account.

The most obvious pattern we see is that support for groups clusters coherently within provinces, with support for al-Qa'ida and the militants fighting in Kashmir clustered in Gujranwala in northeastern Punjab. The second most obvious pattern is that support for the Afghan Taliban is highest in the parts of the country furthest from Afghanistan. Put slightly differently, it appears that, on average, Pakistanis most removed from the consequences of the war in Afghanistan are most supportive of the nonstate actors fighting that war. Critically, Pakistanis in the areas where senior Taliban leadership are thought to reside (Quetta and Zhob) are least tolerant of their organization. Finally, we see that support for the *firqavarana tanzeem* is spread evenly throughout Punjab, not concentrated in Bahawalpur as some have argued.

More specifically, our approach reveals several substantively important patterns. First and foremost, there is clear variation in the intensity of support across divisions within provinces. In the Punjab where respondents are generally weakly positive toward most groups, the intraprovincial variation does not match what one would expect from following ongoing policy discussions. Strikingly, we only find a significant positive endorsement effects for one group, the *firqavarana tanzeems*, in Bahawalpur, which has been widely cited as the locus of militancy within Punjab.¹⁴ If anything Bahawalpur has the lowest average levels of support in this province. Instead, we find strong positive endorsement effects for three of four groups in Gujranwala which has received much less attention from policy makers.¹⁵ We also see positive endorsement effects for the *firqavarana tanzeem* in four of seven divisions within Punjab which is consistent with the idea that many of the sectarian groups draw their support from the Punjab.

Second, areas of Pakistan that suffered a great deal of militant violence before April 2009 show the most negative endorsement effects. Hazara, Kohat, Nasirabad, and Peshawar all suffered particularly heavily from militant violence in early 2009. The results for Quetta are particularly interesting in that this region has long been a safe haven for Afghan Taliban leaders but also suffered heavily from sectarian violence in early 2009. As one might therefore expect, groups espousing explicitly sectarian political goals (al-Qa'ida and the *firqavarana tanzeem*) to fare poorly in public opinion in Quetta. The policy point here is a subtle one insofar as levels of militant activity should not be taken to indicate popular support. Indeed, our analysis suggests that quite the opposite may be true.

¹³For a broader discussion of groups' goals and agendas, see Blair et al. (2011).

¹⁴Bahawalpur district, the center of the division, was listed as one of the 26 USAID focus districts under a now-defunct plan to direct aid to the areas of Pakistan posing the greatest threat of militancy.

¹⁵For example, none of the districts in Gujranwala division were among the 26 USAID focus districts.

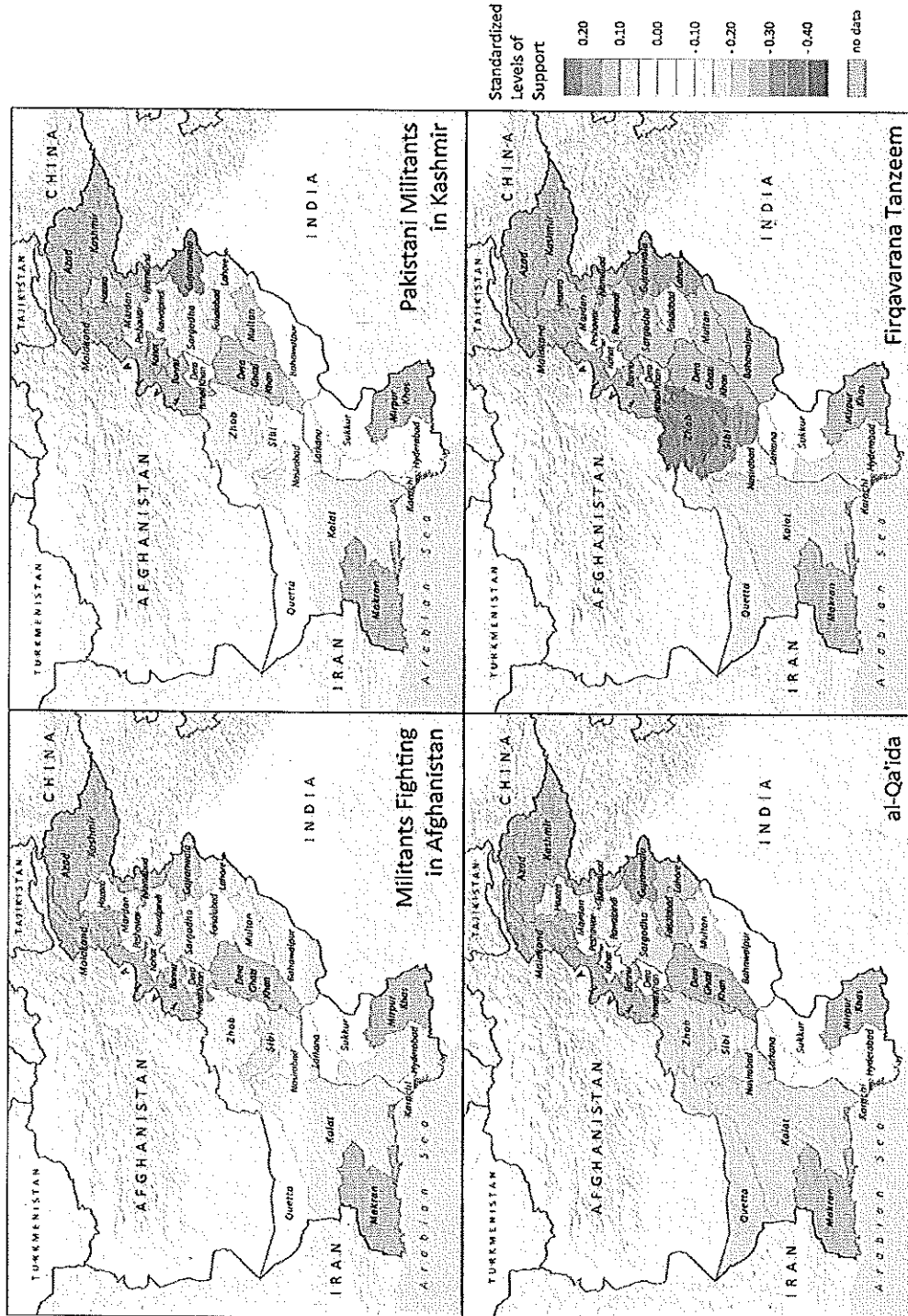


Fig. 4 Map of estimated standardized levels of support for four militant groups in Pakistan. The red represents the estimated positive support, whereas the blue represents the estimated negative support. The dark grey represents the areas where the survey was not conducted and therefore no estimates are available. The estimated levels of support are standardized using the standard deviation of the (posterior) ideal points distribution as a unit.

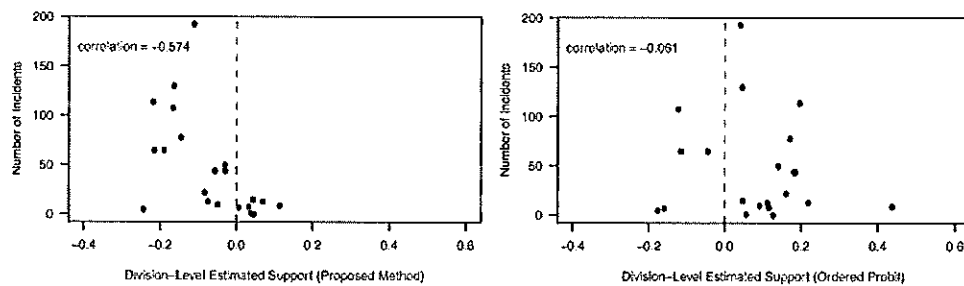


Fig. 5 Bivariate relationship between levels of support for militant groups in Pakistan and incidents of violence within divisions. In the first plot, the horizontal axis represents the estimated level of standardized support within the division, averaged across the four militant groups, obtained from the proposed method. In the second plot, the horizontal axis represents the analogous estimate obtained from using a standard ordered probit model. The vertical axis represents the total number of terrorist incidents in each division from March 2008 through March 2009 (12 months prior to the implementation of the survey) according to the Worldwide Incident Tracking System (WITS) data. There exists a strong negative correlation between the total violence in the area and the estimated average level of support from our model. This correlation is much weaker when average level of support is estimated using the standard ordered probit model.

To probe this insight further, we combine data on attitudes from the survey with panel data on militant violence. The panel data were constructed by georeferencing the 5601 terrorist incidents recorded in Pakistan since 2004 by the National Counterterrorism Center's Worldwide Incident Tracking System (WITS).¹⁶ The WITS data record the location and consequences of all attacks that involve "...premeditated, politically motivated violence perpetrated against noncombatant targets by subnational groups or clandestine agent" but do not list the party responsible for the vast majority of incidents. We examine the bivariate relationship between the estimated support levels and the total number of terrorist incidents from March 2008 through March 2009 (12 months prior to the implementation of the survey).

The first plot of Fig. 5 shows that the impression one takes away from looking at the spatial distribution of violence holds in the data. There is a clear negative relationship between total violence in the year leading up to our survey (vertical axis) and the average level of support for all four groups (horizontal axis) where the correlation is -0.567 with the p value of .007. Although our data do not identify the group that is responsible for each attack, this relationship is strongest for militants focused on Kashmir and the Afghan Taliban, and the correlations between total incidents and support are strongly negative for all four groups.¹⁷ Strikingly, no division that experienced more than 50 incidents in the year before the survey was fielded shows positive support for any militant group.

How do the results based on our proposed model differ from the standard statistical methods that are often used to analyze endorsement experiments? While the direct comparison is somewhat difficult because two models are completely different, we investigate the substantive differences in the resulting estimates of support levels for militant groups. As a standard method, we use the ordered probit model where the responses are regressed on the policy indicator variables, the division indicator variables, the treatment variables, and the interaction of the division indicator variables with the treatment variables, with standard errors clustered by respondent. This simple model assumes that the coefficients for the treatment variables are identical across policy questions and that threshold parameters of response categories are equally spaced across policy questions. We then interpret the coefficients for the treatment variables as scalar measures of support level for militant groups.

The second plot of Fig. 5 shows the relationship between the same terrorist incidence data and the estimated coefficients from the ordered probit model. In contrast to the estimates based on our model, the correlation is weak between total violence in an area and estimated average support for all four militant groups based on a standard ordered probit (-0.129 with p value of .578). Moreover, for each militant

¹⁶We thank Christine Fair and Hamzah Saif for providing this data set.

¹⁷The correlations are: -0.57 (p value = .006) for the Pakistani militant groups in Kashmir, -0.59 (p value = .005) for the Afghan Taliban, -0.47 (p value = .032) for al-Qa'ida, and -0.41 (p value = .062) for the *firqavarana tanzeem*.

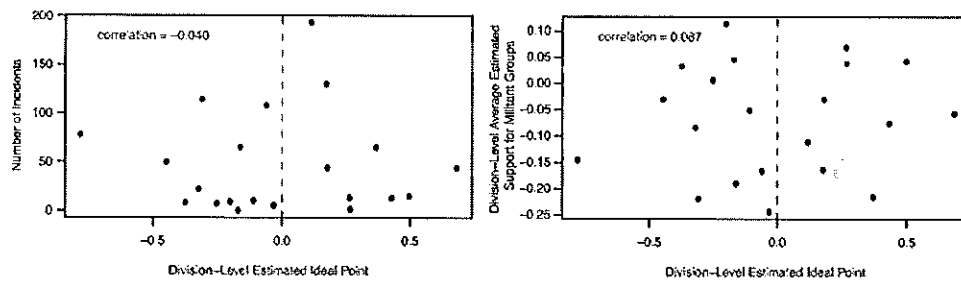


Fig. 6 Bivariate relationship between ideology, incidents of violence and average levels of support for militant groups in Pakistan within divisions. In the both plots, the horizontal axis represents the estimated ideology standardized within the division. In the first column, the vertical axis represents the total number of terrorist incidents in each division from March 2008 through March 2009 (12 months prior to the implementation of the survey) according to the Worldwide Incident Tracking System (WITS) data. In the second plot, the vertical axis represents the estimated level of standardized support within the division, averaged across the four militant groups. From both plots, we see no relationship between division ideology and either division violence or division support for the militant groups.

group, the correlation is considerably lower, and none reaches standard levels of statistical significance.¹⁸ Thus, analysis using a standard ordered probit model would miss a substantively critical result discovered by our model.

Finally, one concern with these results is that ideological differences across divisions may explain the correlation between terrorist incidents and support levels for militant groups. For example, if ideologically liberal areas experience more terrorist attacks and residents of those regions give less support for militant groups, then the negative correlation between violence and support we found above is spurious. To address this possibility, we examine how division-level ideal points in policy space relate to both total violence in an area and to average support for the militant groups in that area. As seen in Fig. 6, there exists no clear relationship between division-level ideal points and total violence. Similarly, there is no strong correlation between ideal points and support at the division level (the corresponding individual-level correlation is also weak, ranging between -0.08 and -0.10). While this in itself is an important substantive finding, the key value of estimating ideal points for this particular analysis is that they allow us to conduct diagnostic checks, showing core empirical results are unlikely to be an artifact of the correlation between violence and ideal points over the underlying policies.

Overall, these results provide evidence for the argument that support for Islamist militancy in Pakistan is concentrated in divisions where there is little violence and in the most economically developed and politically powerful part of the country such as northern Punjab.

5 Simulation Studies

We assess the statistical properties of our model as implemented through JAGS under two circumstances. The first set of simulations is based on the data generating processes that closely follow the observed data of our Pakistan survey experiment. This allows us to examine the reliability of our estimates given in Section 4. The next set of simulations is based upon sample sizes and data generating processes that are more typical among survey experiments in political science. This provides a way to examine whether our proposed method is more broadly applicable. In both sets of simulations, we investigate the frequentist statistical properties of the Bayesian hierarchical measurement models used in our analysis of the Pakistani data, that is, the “Division Model” and the “Division Model with Individual Covariates” as specified in Sections 4.1 and 4.2, respectively.

Below, we report the bias of point estimates (based on posterior mean), the coverage of 90% Bayesian confidence intervals, and the statistical power of $\alpha = 0.1$ level Bayesian hypothesis test. These statistics are computed over 100 simulations. For each simulation run, the model is fitted with three parallel chains

¹⁸The correlations are -0.061 for Pakistani militant groups in Kashmir, -0.365 (p value = .793) for Militants fighting in Afghanistan, 0.021 (p value = .104) for al-Qai'da, and -0.166 (p value = .927) for Firqavarana Tanzeems (p value = .473).

and random starting values. The number of iterations is roughly 60,000. Convergence was diagnosed by two criteria: the values of the \hat{R} statistics and visual inspection of trace and autocorrelation function plots. For expository purpose, we do not provide the detailed information about our simulations (e.g., the true values of the model parameters). Such information and all other details about the data generating process are made available in the replication archive (Bullock, Imai, and Shapiro 2011).

5.1 Monte Carlo Evidence Based on the Pakistan Data

The first set of simulations is based upon the data generating processes that closely follow the observed data from our Pakistan survey experiment. The sample size is 5212, which is larger than typical political science survey experiments. The non-issue-specific parameters (i.e., δ_{division} , $\lambda_{k,\text{division}}$, ω_k , δ^Z , λ_k^Z) are kept constant across all simulations. Throughout the simulations, we use the same distribution of covariates and sample sizes as the Pakistan data.

For each simulation, we sample respondents' ideal points and potential endorsement impacts for all groups and policies according to our model (conditional on their covariates when appropriate). After splitting the sample into equal halves control group and treatment group, the sequence of endorsers is randomized for the treatment group. We then generate the responses to a set of questions, each of which has a random set of cut-points and discrimination parameter according to our model specification. In addition, the true values of the model parameters are chosen such that these responses were to fall within a plausible range. This constitutes the generic data generating process.

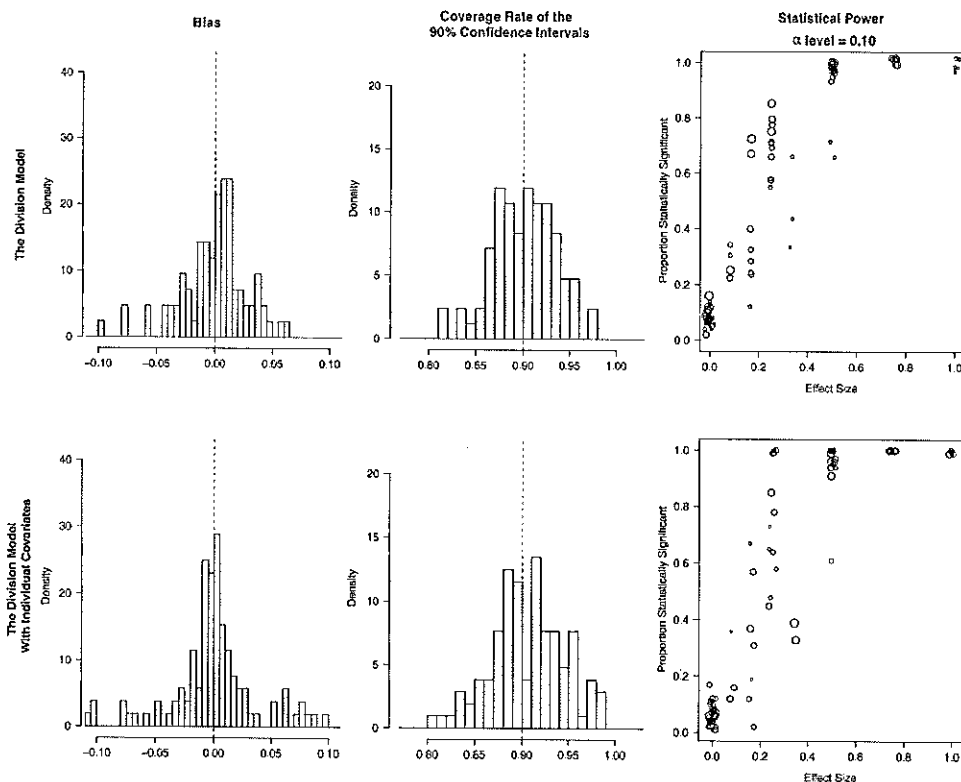


Fig. 7 Simulation results based on the Pakistan data. Across both models, the distribution of bias across parameters (left column) centered around zero is small in magnitude relative to the standard deviation of respondents' ideal points. The coverage rates of the 90% Bayesian confidence intervals across parameters (middle column) are also centered around the nominal rate. Finally, the statistical power of the proposed models (represented as the proportion of statistically significant estimates in the right column) is quite high and increase rapidly with effect size. Moreover, we observe that the statistical power is a function of sample size within each division, which is reflected by the size of circles.

As with any set of simulations testing a statistical model, we would ideally observe three criteria. First, there should be no bias in the recovered parameter estimates. Second, the recovered confidence intervals should be of appropriate width. Confidence intervals that are too short will lead us to have inappropriately high confidence in our point estimates, whereas intervals that are too long will leave us unable to make clear distinctions from the data. Finally, given the traditional focus on hypothesis testing, we want our tests to have large statistical power if possible. That is, we would like to detect the departure from the null hypothesis whenever it exists.

The results from this set of simulation studies are summarized in Fig. 7. For the Division Model, we observed an average bias of -0.007 and an average root mean squared error (RMSE) of 0.138 . The coverage rate of the 90% confidence intervals was 0.887 . Finally, the statistical power of the proposed models (represented as the proportion of statistically significant at the $\alpha = 0.1$ level in the right column) is quite high and increases rapidly with effect size and division level sample size (represented by the size of circles in the top right corner graph). Moreover, the Type I (false positive) error rate was 0.099 , close to the nominal level of $\alpha = 0.1$.

The results for the Division Model with Individual Covariates are equally impressive in terms of all three criteria, except that the statistical power may be low for some parameters. Overall, these models seem to recover the true parameter values efficiently and their confidence intervals have appropriate coverage rates.

5.2 Monte Carlo Evidence with Varying Sample Sizes

Most political science data sets are much smaller than the Pakistan data set, and more often, the responses to survey questions have central response distributions. To give a better understanding of the performance of our model in typical situations, we conduct simulation studies with varying sample sizes. In particular, we repeat the data generating process as above, but with three exceptions. First, the number of geographic

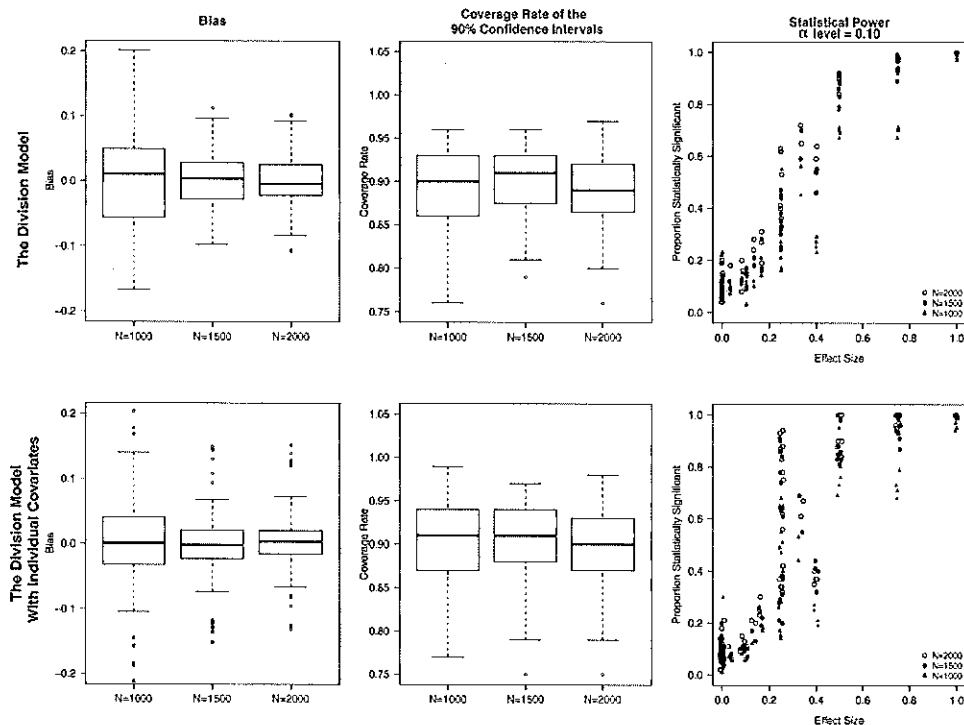


Fig. 8 Simulation results with varying sample sizes. Each column represents a criteria by which the simulation results should be judged, while each row represents a different model. As with the previous set of simulations, we find little evidence of bias or inadequate coverage. As expected, however, the statistical power does decrease as the sample size gets smaller.

Table 1 Simulation summary statistics. For each set of simulations, the average bias across all parameters of interest is negligible. Upon first glance, the average RMSE is not as small as one might hope. But the coverage statistics of the 90% confidence intervals are of appropriate size, which suggests that the errors are not necessarily problematic. Moreover, the Type I (false positive) error rate is close to the nominal level of $\alpha = 0.1$

Model	<i>N</i>	<i>Average bias</i>	<i>Average RMSE</i>	<i>Average coverage</i>	<i>Type I error</i>
Division	5212	-0.001	0.137	0.907	0.063
Division with individual covariates	5212	0.000	0.123	0.911	0.066
Division	1000	0.006	0.251	0.888	0.081
Division	1500	0.003	0.191	0.900	0.080
Division	2000	-0.003	0.178	0.890	0.086
Division with individual covariates	1000	-0.001	0.224	0.900	0.069
Division with individual covariates	1500	-0.003	0.173	0.901	0.071
Division with individual covariates	2000	0.004	0.163	0.894	0.071

units was reduced from 26 to 16 divisions, with equal sample sizes across divisions. Second, the sample size is reduced considerably to sizes 1000, 1500, and 2000. Third, the distribution of responses has a stronger central tendency on the 1–5 scale.

Figure 8 presents the simulation results based on all three models. These models perform reasonably well with small bias (left column) and appropriate coverage probabilities (middle column). In particular, as the sample size increases, the bias gets smaller and the coverage rate improves. As expected, the statistical power (right column) does decrease as the sample size gets smaller. However, even when the sample size is as small as 1000, the statistical power is reasonable; we reject the null hypothesis of zero effect about half of the time for the parameter whose true value is around 0.4.

5.3 Summary of Simulation Results

Our ability to recover precise unbiased parameter estimates gives us confidence that our estimation procedure is well suited for endorsement experiments. Although no set of simulations can cover all plausible situations a researcher might face, our model performs well across a broad range of scenarios. Table 1 shows that the proposed models have good frequentist properties on several criteria. For each set of simulations, the average bias (first column) across all parameters of interest is negligible. The average RMSE (second column) is larger than one might hope. But the coverage statistics of the 90% confidence intervals (third column) are of appropriate size, which suggests that the errors are not necessarily problematic. Moreover, the Type I (false positive) error rate is consistently close to the nominal level of $\alpha = 0.1$ (last column). When compared to ordinal logistic regression and other standard estimation techniques, our model performs well. For example, using ROC curves to assess in-sample and out-of-sample prediction rates, our model is equally or slightly more efficient while remaining unbiased (data not shown).

6 Concluding Remarks and Suggestions for Applied Researchers

Endorsement experiments are a survey methodology designed to ascertain levels of political support for socially sensitive actors. While their effectiveness must be empirically validated more extensively in the future, endorsement experiments may provide researchers with one way to indirectly ask sensitive questions, thereby avoiding the problems of social desirability bias and item nonresponse. The Pakistan experiment analyzed in this paper is such an instance where directly asking about support for militant groups in many parts of the country would have raised safety concerns for survey respondents and enumerators as well as the possibility of bias and high nonresponse.

In this paper, we develop a new Bayesian hierarchical measurement model to analyze endorsement experiments by appropriately combining multiple policy questions and yielding estimates of support within the item response theory framework so that they are interpretable. Through our empirical application, we show that our Bayesian measurement model can effectively recover subtle regional differences in support levels, yielding new substantive results which standard methods fail to find. The simulation studies

demonstrate that across different sample sizes, the proposed methodology has good frequentist properties: low bias, confidence intervals with good coverage, and relatively large statistical power.

Applying our model to the data from Pakistan revealed a number of key facts that were not apparent from the commonly used statistical models. We showed that attitudes toward militant groups are highly regionally clustered and that support is generally lower where there have been a large number of militant attacks. Future research should examine how violent incidents affect the support levels for (perceived) perpetrators by incorporating this information directly into the proposed model. In addition, we find that individual-level variables that are the focus of much policy attention—many of the aid programs directed at addressing the threat from Islamist militancy focus on increasing education and income, for example—matter little once regional differences in support levels are taken into account.

Despite these advantages, it is important to recognize the limitations of endorsement experiments and the proposed methodology. In particular, as discussed in Section 3.2, inferring level of political support from endorsement experiments requires the critical assumption that difference in respondents' support for policies induced by endorsement can be attributed to support for the group rather than respondents' different interpretations of policy questions (e.g., learning about underlying policies). Since the observed data alone cannot verify this assumption, we provide the following practical recommendations for applied researchers who are planning to use endorsement experiments in their research.

- Ask survey respondents to locate each policy before assessing support and examine whether endorsements affect perceived policy locations. If it does, that would suggest that respondents in the treatment group are learning about policies, making it difficult to infer support level from endorsement experiments.
- Choose policies that are well known among respondents and are on the same one-dimensional policy space.¹⁹ This minimizes the possibility that endorsement provides new information about these policies and avoids the statistical model with a high-dimensional policy space that is difficult to estimate and interpret.²⁰
- Measure the level of political knowledge each respondent has about the selected policies. Such a measure allows researchers to focus their analysis on the subset of respondents who are less likely to learn about policies themselves from endorsement.

Using surveys to measure political support for socially sensitive actors creates many challenges with respect to the design, implementation, analysis, and interpretation of surveys. While there is no perfect solution, endorsement experiments provide a possible way to reduce bias resulting from social desirability and nonresponse. We show that it is possible to analyze endorsement experiments within the item response theory framework and obtain measures of support levels on the same scale as ideal points of survey respondents. The proposed model efficiently combines responses across different policy questions and shows how much the endorsement by each actor can shift the ideal points of respondents on the common ideological dimension. Equipped with the new tools, we hope that researchers can now effectively apply endorsement experiments to measure political support when direct questioning is either impossible or unlikely to yield reliable measures.

Funding

Financial support from the National Science Foundation grant (SES-0849715 for Imai), Department Homeland Security (2007-ST-061-000001 for Shapiro) through the National Center for Risk and Economic Analysis of Terrorism Events, and the Mamdouha S. Bobst Center for Peace and Justice (for Imai and Shapiro) is acknowledged. The Institute for Quantitative Social Science at Harvard University provided the computational support necessary for simulation studies. We thank Simon Jackman, John

¹⁹Note that in the literature on source cues, researchers often choose unknown policies in order to examine how much respondents' opinion can be influenced by endorsements (see e.g., references in footnote 1). These researchers are not necessarily interested in interpreting these changes as support for endorsers.

²⁰The question of how the performance of the proposed statistical model is affected by these deviations from its assumptions is left to future research.

Londregan, Nolan McCarty, Stephen Nicholson, Michael Peress, and seminar participants at California Institute of Technology, Columbia, New York University, and Princeton for helpful comments. Christine Fair and Neil Malhotra helped design the survey, to which we apply the statistical model presented in this paper. We also thank Joshua Borkowski for his help in creating the map presented in Fig. 4.

References

- Bafumi, J., A. Gelman, D. K. Park, and N. Kaplan. 2005. Practical issues in implementing and understanding Bayesian ideal point estimation. *Political Analysis* 13:171–87.
- Bafumi, J., and M. Herron. 2010. Leapfrog representation and extremism: A study of American voters and their members in Congress. *American Political Science Review* 104(3):519–42.
- Blair, G., C. C. Fair, N. Malhotra, and J. N. Shapiro. 2011. Poverty and support for militant politics: Evidence from Pakistan. SSRN Working Paper.
- Blair, G., and K. Imai. 2010. Statistical analysis of list experiments. Working paper. <http://imai.princeton.edu/research/listP.html>.
- Bullock, W., K. Imai, and J. N. Shapiro. 2011. Replication data for: Statistical analysis of endorsement experiments: Measuring support for militant groups in Pakistan. hdl:1902.1/14840. The Dataverse Network.
- Clinton, J., S. Jackman, and D. Rivers. 2004. The statistical analysis of roll call data. *American Political Science Review* 98:355–70.
- Clinton, J. D., and D. E. Lewis. 2008. Expert opinion, agency characteristics, and agency preferences. *Political Analysis* 16:3–20.
- Cohen, G. L. 2003. Party over policy: The dominating impact of group influence on political beliefs. *Journal of Personality and Social Psychology* 85:808–22.
- Fair, C. C., N. Malhotra, and J. N. Shapiro. 2009. The roots of militancy: Explaining support for political violence in Pakistan. Working Paper. Princeton University.
- Gelman, A., and J. Hill. 2007. *Data Analysis Using Regression and Multilevel/Hierarchical Models*. Cambridge: Cambridge University Press.
- Gelman, A., and D. B. Rubin. 1992. Inference from iterative simulations using multiple sequences (with discussion). *Statistical Science* 7:457–72.
- Gingerich, D. W. 2010. Understanding off-the-books politics: Conducting inference on the determinants of sensitive behavior with randomized response surveys. *Political Analysis* 18:349–80.
- Heckman, J. J., and J. M. Snyder. 1997. Linear probability models of the demand for attributes with an empirical application to estimating the preferences of legislators. *RAND Journal of Economics* 28:142–89.
- Holland, P. W. 1986. Statistics and causal inference (with discussion). *Journal of the American Statistical Association* 81:945–60.
- Imai, K. 2011. Multivariate regression analysis for the item count technique. *Journal of the American Statistical Association* 106:407–16.
- Kam, C. D. 2005. Who toes the party line? Cues, values, and individual differences. *Political Behavior* 27:163–82.
- Lax, J. R., and J. H. Phillips. 2009. How should we estimate public opinion in the states? *American Journal of Political Science* 53:107–21.
- Muraki, E. 1992. A generalized partial credit model: Application of an EM algorithm. *Applied Psychological Measurement* 16:159–76.
- Nicholson, S. P. Forthcoming 2011. Dominating cues and the limits of elite influence. *Journal of Politics*.
- Pakistan Federal Bureau of Statistics. 2008. Labour force survey 2007–8. Technical report.
- Park, D. K., A. Gelman, and J. Bafumi. 2004. Bayesian multilevel estimation with poststratification: State-level estimates from national polls. *Political Analysis* 12:375–85.
- Peress, M., and A. Spirling. 2010. Scaling the critics: Uncovering the latent dimensions of movie criticism with an item response approach. *Journal of the American Statistical Association* 105:71–83.
- Plummer, M. 2009. *JAGS: Just Another Gibbs Sampler*. <https://sourceforge.net/projects/mcmc-jags>.
- Poole, K. T., and H. Rosenthal. 1985. A spatial model for legislative roll call analysis. *American Journal of Political Science* 29:357–84.
- Rivers, D. 2003. Identification of multidimensional spatial voting models. Unpublished manuscript, Department of Political Science, Stanford University.
- Shapiro, J. N., and C. C. Fair. 2010. Why support Islamist militancy? Evidence from Pakistan. *International Security* 34:79–118.
- Sniderman, P. M., and T. Piazza. 1993. *The Scar of Race*. Cambridge: Harvard University Press.
- United States Agency for International Development. 2009. Development assistance and counter-extremism: A guide to programming. Technical report. http://www.usaid.gov/locations/sub-saharan.africa/publications/docs/da_and_cea_guide_to_programming.pdf.