

Using Hierarchical Models to Estimate Heterogeneous Effects

Joshua Alley
Assistant Professor
University College Dublin
joshua.alley@ucd.ie

May 22, 2024

Abstract

This note describes why, when, and how to use Bayesian hierarchical models to estimate heterogeneous effects. While an ample literature suggests that hierarchical models provide helpful regularization and information about variation, political scientists rarely use them to estimate heterogeneous effects. Doing so is simple, however. To start, specify groups based on quantities of interest such as demographics, context, and policy relevance. Then, fit a hierarchical model where treatment slopes and intercepts vary across groups. This captures systematic and random variation in heterogeneous effects, estimates effects within each group, and measures effect variance. Hierarchical modeling provides an intermediate tool between interactions or subgroup analyses and machine learning approaches to discovering complex heterogeneity. It is more flexible than interactions and reduces the risk of underpowered subgroup comparisons. At the same time, it is more theoretically informed and interpretable than some machine learning approaches, as well as easier to implement in small datasets. I document these claims with a simulation analysis and extension of a published study. Researchers can thus use hierarchical models alongside other approaches to understand heterogeneous effects for scholarship and policy.

1 Introduction

Whether in observational or experimental studies, every independent variable social scientists examine impacts some units differently than others. Common estimands aggregate heterogeneous effects.¹ These average effects are useful, but they often obscure interesting and important variation.

As a result, understanding heterogeneous effects is essential for policy and scholarship. Estimating heterogeneity allows scholars to clarify the connection between their independent variable and outcome. Policymakers can maximize the impact of finite resources with targeted interventions, for example by providing job training to individuals who are more likely to benefit.

This letter explains why, when and how to use hierarchical models to estimate heterogeneous effects. I identify when researchers can profitably use hierarchical models, and when other tools make more sense. A large statistics literature suggests that Bayesian hierarchical models are a useful tool for heterogeneous effects estimation (e.g., Feller and Gelman (2015); McElreath (2016); Dorie et al. (2022)). Political scientists tend to rely on interactions or machine learning tools instead, however. For instance, of the three applied political science citations of Feller and Gelman (2015), only Marquardt (2022) models treatment effects.

The main advantage of hierarchical modeling is regularization, as partial pooling pulls estimates towards an overall mean. This matters because many political science applications have low power even for main effects (Arel-Bundock et al., 2022). Adequately powered estimates of interactions sometimes requires even more data (Gelman, 2018), which is often unavailable. As a result, statistically significant heterogeneous effect estimates may be far too large—the result of noise in the data, not systematic differences. This is one reason why findings based on interactions are especially unlikely to replicate (Simmons, Nelson and Simonsohn, 2011).

¹For instance, Abramson, Koçak and Magazinnik (2022) note that the average marginal component effect (AMCE) of conjoint experiments gives more weight to intense preferences.

At the same time, social scientists often posit conditional theories. The resulting accumulation of single-modifier arguments suggests that multiple factors generate heterogeneous effects. For example, scholarship on audience costs has considered how foreign policy dispositions (Kertzer and Brutger, 2016), partisanship (Levendusky and Horowitz, 2012), gender (Barnhart et al., 2020; Schwartz and Blair, 2020) and policy preferences (Chaudoin, 2014) all modify individual reactions to a leader backing down from a threat. Multiple modifiers further exacerbate power concerns.

Hierarchical modeling of heterogeneous effects addresses these concerns, and thereby fills a gap between interactions and machine learning.² Parametric interactions and subgroup analyses are ubiquitous because they are easy to interpret, but these are subject to power concerns. More recent work employs random forests (Green and Kern, 2012; Wager and Athey, 2018), support vector machines (Imai and Ratkovic, 2013), and ensemble methods (Grimmer, Messing and Westwood, 2017; Künzel et al., 2019; Dorie et al., 2022). These machine learning algorithms often have some regularization component and can discover complex patterns and high-dimensional variation, but can be difficult to interpret and implement, especially in smaller social science datasets.

Using a hierarchical model is more flexible than parametric interactions but easier to implement and interpret than machine learning approaches. It preserves a simple and interpretable structure, while accommodating more factors and ameliorating the downsides of subgroup analysis via partial pooling. This facilitates argument testing. Unlike machine learning, the hierarchical approach lacks the flexibility to discover high-dimensional heterogeneity, however.

Hierarchical modeling therefore works best when there are more than two modifying factors and therefore many subgroups of interest, as well as less emphasis on discovery. These

²Blackwell and Olson (2022) describe a lasso approach to interactions that sits between machine learning and linear regressions.

models are best at capturing variation across groups and levels when there are multiple potential modifiers. This also works well when researchers have a clear sense of the relevant groups.³

There are two key steps when theory and data make using hierarchical models worthwhile. First, researchers should define groups based on potential sources of heterogeneity such as other treatments, context, demographics, or policy concerns. Second, they should estimate heterogeneous effects across those groups using a hierarchical model with varying slopes and intercepts for every unique group. Modeling heterogeneous effects in this way produces interpretable results, which facilitates argument testing. It also allows researchers to examine effects within groups, compare different sources of heterogeneous effects and describe how much an effect varies.

While frequentist estimation of hierarchical models is possible, Bayesian estimation is straightforward and more informative. Bayesian estimation provides crucial information by connecting parameters through common prior distributions, thereby regularizing estimates and propagating uncertainty. Working with posterior distributions also gives researchers more flexibility to present diverse information about how and when effects vary. While computation and coding were once a barrier to employing Bayesian methods, fitting a wide range of hierarchical models is straightforward with the `brms` package in `R` (Bürkner, 2017).⁴ Calculating substantive effects is also simple (Arel-Bundock, N.d.).

In the remainder of this note, I describe how and when to estimate hierarchical models of heterogeneous effects and demonstrate the process by analyzing a study of how military alliances shape public support for war by Tomz and Weeks (2021). The reanalysis reveals that alliances increase support for intervention most among white men who support international engagement but are otherwise skeptical of using force.

³(Goplerud, 2021) introduces a model that uses Bayesian structured sparsity to estimate which group coefficients are similar and which are different. In this researchers use theory to inform the potential sets of groups, but common estimates for groups are data driven.

⁴I provide example code in this note and the appendix.

2 Hierarchical Modeling of Heterogeneous Effects

There are two steps in hierarchical models of heterogeneous effects. First, researchers must define the groups over which an independent variable's impact changes. Unique combinations of characteristics such as other treatments, context and demographics determine groups.

Researchers should create groups based on what variation is most important and interesting. Theory, policy concerns, or normative factors are all possible motivations.

Setting groups is the most important task, because it determines what heterogeneous effects a researcher estimates. Defining groups before model fitting defines what variation is most important, links heterogeneous effects to theory, and structures modeling.⁵ Defining groups without careful thought risks obfuscating results and can hinder model fitting.

There are three general approaches to defining groups. First, researchers can set groups using combinations of other treatments, especially when an intervention has several dimensions but theory emphasizes one of them. The experimental design determines groups, and the model estimates heterogeneous treatment effects. If researchers want to know how different issues shape the impact of elite foreign policy cues (Guisinger and Saunders, 2017), they could define groups by issues, for instance. Hierarchical estimators for topic-sampling experiments apply this idea to estimate how a treatment varies across different topics (Clifford and Rainey, 2023). Researchers sometime use fully crossed interactions to estimate the impact of a treatment across experimental strata, but this approach risks spurious results due to inadequate power.

A second approach uses unit, demographic and contextual factors to create groups and estimate effect heterogeneity. Here, researchers examine what factors within or around units shape their response to an independent variable. For example, Alley (2021) uses alliance characteristics to examine when alliance membership increases or decreases military spending. Other use cases include estimating how different demographic groups or geographic units respond

⁵It also facilitates pre-registration when applicable.

to an intervention. One might examine how the impact of a national-level intervention varied across states, for example.⁶

Third, researchers might use hierarchical models to address specific policy concerns. Policy analysts often want to know how an intervention impacts a specific population. Researchers might want to know if a job-training program improves employment outcomes for black women in the South, for instance.

Whether researchers use other treatments, context, or policy to determine groups, the number of grouping factors depends first on theory. There are some practical constraints, however. Using too many factors can lead to model fitting and interpretation problems by creating many small groups. How many factors is too many depends especially on the data—some datasets can support reasonably large groups for many factors. At the other extreme, using only one grouping factor will add relatively little value compared to interactions.

After defining groups, the second step is fitting a hierarchical model of effects within groups. The model employs a mix of varying intercepts and slopes to estimate the impact of an intervention in each group.⁷

This approach can address diverse problems, but for ease of exposition consider making between-unit comparisons based on an experimental treatment. Start with N units indexed by i , some of which receive a binary treatment T . Assume that the outcome variable y is normally distributed with mean μ_i and standard deviation σ .⁸ g indexes the researcher-defined groups, which include two variables $\nu 1$ and $\nu 2$ and their interaction.

The outcome for each unit is then a function of group varying intercepts α_g , an optional matrix of control variables \mathbf{X} , and a set of group treatment effects θ_g , which are normally

⁶Extrapolation to a representative sample for such units might require poststratification.

⁷If other units such as states define the groups, rather than combinations of modifying variables, then adding group-level predictors is essential. For example, in a model where an effect varies by state, adding state-level variables like ideology, population and GDP would avoid partially pooling small groups too far towards the overall mean.

⁸Researchers should use binary, categorical and other outcome likelihoods.

distributed with mean η_g and standard deviation σ_θ . The researcher divides all units into g groups based on unique combinations of heterogeneous effect predictors \mathbf{Z} . Each θ parameter estimates how the treatment effect varies across the values of each variable. To capture correlations between the random intercepts and varying slopes $\rho\sigma_\alpha\sigma_\theta$, these variables should have a common multivariate normal prior.

$$y_i \sim N(\mu_i, \sigma) \quad (\text{Likelihood})$$

$$\mu_i = \alpha + \alpha_{v1} + \theta_{v1}T + \alpha_{v2} + \theta_{v2}T + \alpha_{v1*v2} + \theta_{v1*v2}T + \mathbf{X}\beta \quad (\text{Outcome Equation})$$

$$\begin{pmatrix} \alpha_{v1} \\ \theta_{v1} \end{pmatrix} \sim N \left[\begin{pmatrix} \mu_{v1}^\alpha \\ \mu_{v1}^\theta \end{pmatrix}, \begin{pmatrix} \sigma_\alpha^2 & \rho\sigma_{v1}^\alpha\sigma_{v1}^\theta \\ \rho\sigma_{v1}^\alpha\sigma_{v1}^\theta & \sigma_\theta^2 \end{pmatrix} \right] \quad (\text{Common Prior: V1})$$

$$\begin{pmatrix} \alpha_{v2} \\ \theta_{v2} \end{pmatrix} \sim N \left[\begin{pmatrix} \mu_{v2}^\alpha \\ \mu_{v2}^\theta \end{pmatrix}, \begin{pmatrix} \sigma_\alpha^2 & \rho\sigma_{v2}^\alpha\sigma_{v2}^\theta \\ \rho\sigma_{v2}^\alpha\sigma_{v2}^\theta & \sigma_\theta^2 \end{pmatrix} \right] \quad (\text{Common Prior: V2})$$

$$\begin{pmatrix} \alpha_{v1*v2} \\ \theta_{v1*v2} \end{pmatrix} \sim N \left[\begin{pmatrix} \mu_{v1*v2}^\alpha \\ \mu_{v1*v2}^\theta \end{pmatrix}, \begin{pmatrix} \sigma_\alpha^2 & \rho\sigma_{v1*v2}^\alpha\sigma_{v1*v2}^\theta \\ \rho\sigma_{v1*v2}^\alpha\sigma_{v1*v2}^\theta & \sigma_\theta^2 \end{pmatrix} \right] \quad (\text{Common Prior: V1*V2})$$

$$\theta_g = \theta_{v1}V1 + \theta_{v2}V2 + \theta_{v1*v2}V1*V2 \quad (\text{Group Slopes})$$

$$\alpha_g = \alpha_{v1} + \alpha_{v2} + \alpha_{v1*v2} \quad (\text{Group Intercepts})$$

This approach lets slopes vary across multiple variables within a group. The net impact of the treatment in each group depends on the linear combination of slopes in that group, specifically the θ parameters and values of each variable. Similarly, the group intercepts can be calculated as the sum of the corresponding random intercepts for each grouping variable.⁹

⁹In brms for a model with no controls and two variables modifying the impact of a treatment, the model formula is $y \sim 1 + (1 + \text{treat} | \text{var1}*\text{var2})$

The above model can be fit with Bayesian or frequentist methods, but Bayesian estimation offers important advantages. First, it is more flexible, and including prior information can facilitate model fitting and convergence. Priors also help regularize estimates by pulling extreme groups towards the overall mean. Working with posterior distributions also provides a wealth of information about effect heterogeneity and propagates uncertainty.

In interpreting these models, researchers should leverage the full range of information from the different parameters. First, the θ posteriors give the impact of a variable within each group. All θ s reflect a systematic component from the predictors in $\mathbf{Z}\lambda$ and a random variation in slopes from σ_θ . Unless the group-level predictors are weak correlates of treatment response, the systematic component will dominate. The random variation is similar to the error term in regression- it expresses how much variation is left in addition to the systematic component.

In addition to group-specific effect estimates, a hierarchical model facilitates rich description of effects across groups. Researchers can calculate variance in the θ parameters across groups and compare the posteriors of different θ s. The σ_θ parameter summarizes the random variation. Other techniques such as interactions in OLS with robust standard errors provide less information.

3 When to Use Hierarchical Models

In deciding whether to use a hierarchical model, researchers must weigh its unique advantages and disadvantages. In general, estimating heterogeneous effects in this way has three advantages. First, researchers can make detailed inferences about heterogeneous effects in an interpretable framework. This helps examine theories that predict how an effect varies and compare sources of variation.¹⁰ Partial pooling also facilitates reasonable estimates for small groups by sharing information across groups and incorporating predictors in the heteroge-

¹⁰Rescaling variables in the heterogeneous effects equation can aid model fitting and coefficient comparisons (Gelman, 2008).

neous effects equation. Finally, this approach will be faster than machine learning approaches for many datasets, easier to use in small datasets, and may scale better than models of individual treatment effects.

Like all methods, the hierarchical approach has downsides, some of which can be ameliorated with modifications, while others should lead researchers to use different tools. Because groups are based on unique combinations of heterogeneous effect variables, using multiple continuous variables in the heterogeneous effects equation creates many small groups or individual treatment effects, which increases the risk of sampling problems, especially in small datasets. If using continuous variables hinders model convergence, researchers can bin continuous variables.

Furthermore, hierarchical models can show general trends, but will not make powerful comparisons between every group. Researchers who want to compare specific groups may lack empirical leverage, especially for small groups.

With these considerations in mind, when should researchers use hierarchical models in place of interactions? If only one factor modifies an effect, interactions are best, as the extra information hierarchical models provide is less valuable.

With two or more modifiers, hierarchical models begin to add value beyond interactions. Interpreting triple interactions between a variable and two modifiers is challenging. The advantages of hierarchical modeling increase with the number of modifiers, until additional modifiers create small groups that complicate model fitting. The thresholds where the number of modifiers becomes an issue for hierarchical modeling depends on the data, as larger datasets can support more groups.

The relative use cases of hierarchical models and machine learning are different. Unlike machine learning approaches, hierarchical models will not discover high-dimensional interactions. Researchers can add flexibility with additional interactions or non-linear specifications in either level of the model, but this requires a priori specification. Therefore, if researchers

	Hierarchical Models	Interactions/Subgroup	Machine Learning
Factors	Two or more	One or two	Many
Sample Size	Conditional on number of factors	Medium to large, depending on main effect size	Large
Complexity	Medium	Low	High
Computational Cost	Medium	Low	High
Interpretability	High	High	Low
Modifiers	Specified	Specified	Discovered or Specified

Table 1. *Key characteristics of different approaches to estimating heterogeneous effects.*

want to focus on flexible discovery, not testing an argument with multiple sources of treatment heterogeneity, they should rely more on machine-learning.

In summary, researchers should continue to use interactions for single modifiers and machine learning to discover complex interactions. Hierarchical modeling works well when there are two or more modifiers and researchers have adequate data to support an informative model. Table 1 summarizes some relevant characteristics of hierarchical, interaction and machine learning approaches to heterogeneous effects. Hierarchical modeling is thus an intermediate tool between interactions and machine-learning, where researchers need more flexibility than interactions but are not willing or able to tackle the computational and interpretation challenges of machine learning.

4 Performance on Simulated Data

To assess the performance of hierarchical models compared to OLS interactions, I first compare their performance on simulated data. Simulating varying effects across multiple groups is challenging. To begin, I define a binary treatment variable as randomly drawn from a binomial distribution with a probability of .5. I then define three randomly generated grouping

variables; one four category multinomial variable with equal probability of drawing each category, a binomial variable with a 70% chance of drawing one, and second binomial variable with a 30% chance of drawing one. There are 16 unique combinations of the values in these groups.

The corresponding 16 group-level effects are a function of fully crossed interactions between the grouping variables and treatment. I simulate the interaction effects by drawing them from a normal distribution with a mean of .15 and standard deviation of either .05, .25 or .75. Increasing the standard deviation adds greater variation to the group-level effects.

Those interactions predict the mean of the outcome variable, μ_y , which has a standard deviation of .25. I vary the sample size of this outcome and the predictors from 1,000, 2,500, or 5,000 observations. This provides more data to compare estimator performance and assess potential speed tradeoffs with a hierarchical model.

In total, the simulation varies both the standard deviation of the interactions and the sample size. This assesses the regularization benefits of varying slopes in situations where more data and greater effect variability may make it more necessary. More variation in effects may make regularization more impactful, while larger datasets will make it less helpful.

5 Example Application: Alliances and Public Support for War

In the following, I further demonstrate how the hierarchical approach works and the benefits of regularizing effect estimates by reanalyzing a study by Tomz and Weeks (2021) (TW hereafter). TW examine whether the public is more willing to go to war for an allied country. In a factorial experiment with vignettes, they find a 33% average increase in support for military intervention on behalf of another country if that country is an ally. This is a large and potentially important relationship.

Given the size of the main effect, TW’s paper is in some ways a best case scenario for comparing interactions and hierarchical models. Corresponding interaction effects may be quite large, and their sample size of 1,200 respondents is not unusual in published work. At the same time, TW estimated an array of interactions to check how other treatments modify the impact of alliances. There are 64 unique treatment groups with anywhere from 11 to 32 respondents in each, so estimates of the impact of alliances in the 32 pairs of alliance treatment and control groups employ at most 54 data points. As such, employing varying slopes for regularization will likely offer substantial benefits. I document these gains first by analyzing the how other experimental treatments modify the impact of alliances, and then exploring demographic differences by alliance.

5.1 Differences by Experimental Scenario

Along with alliances, TW randomly assign whether the potential beneficiary of U.S. intervention is a democracy or not, the stakes of intervention, the potential costs, and the region of the world. They estimate the impact of alliances in the 32 treatment conditions with an OLS model that fully crosses interactions between the treatments, and calculate marginal effects that average over these groups. I keep the same fully crossed structure in the treatment interactions to define groups, but use a varying slopes model to estimate the impact of alliances in each treatment group.

Figure 1 compares the estimated alliance treatment effects across the experimental groups with the OLS and hierarchical models. Two aspects of Figure 1 show the regularization benefits of hierarchical modeling. First, the hierarchical estimates are more precise. The credible intervals in the hierarchical model are smaller because they draw on information from every group.

Second, the hierarchical estimates are less variable. This reduces the estimated variation in how alliances impact mass attitudes, and is obvious if we compare estimates within regions.

Heterogeneous Alliance Treatments Divided By Experimental Conditions

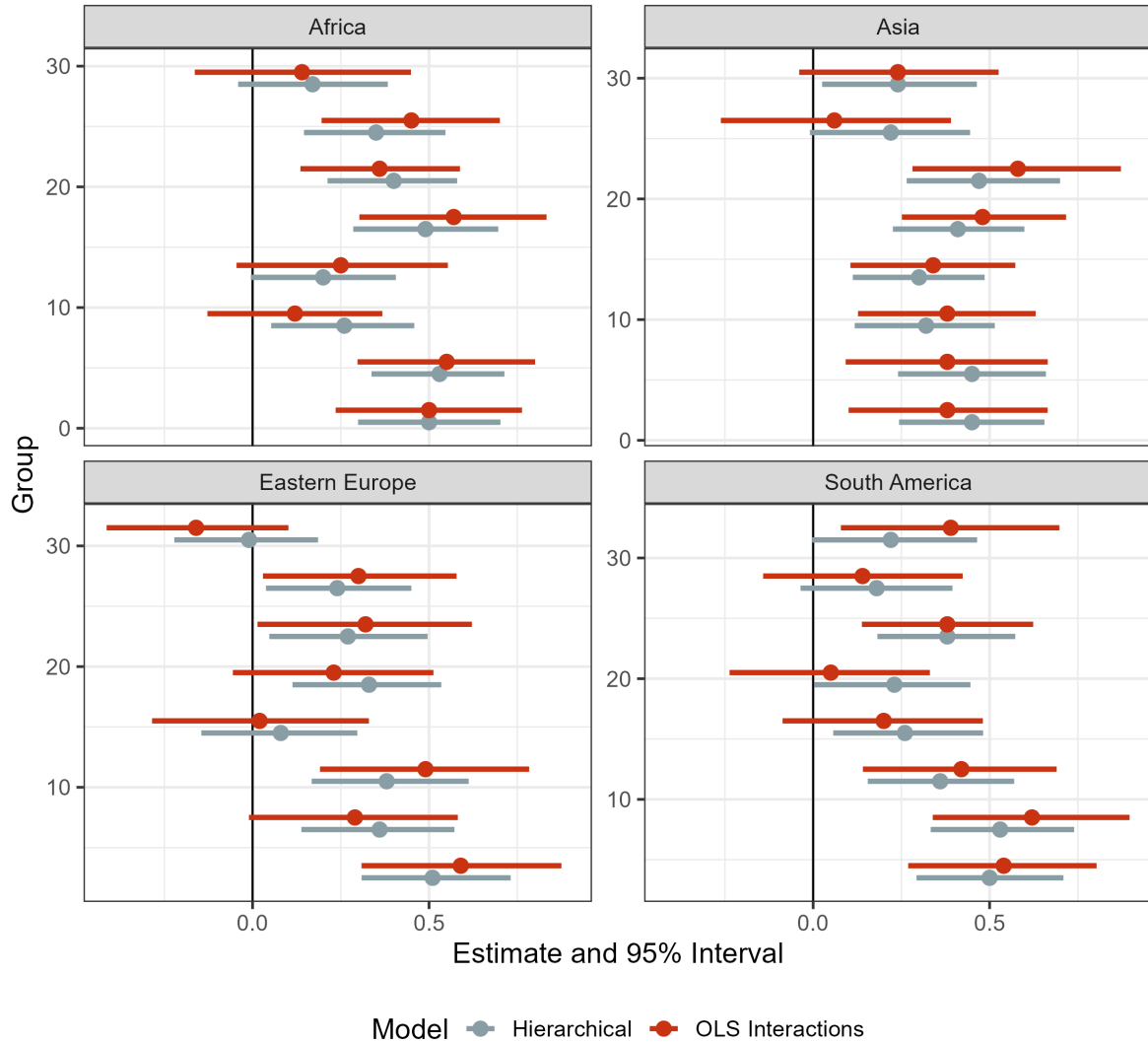


Figure 1. Comparison of OLS and hierarchical estimates of the impact of alliance across experimental conditions. Estimates divided based on the regional treatment variable for ease of presentation.

The OLS interactions are more dispersed, while the hierarchical estimates hew more closely to the overall mean. This is especially notable in the African and Latin American scenarios. Inasmuch as differences across scenarios are driven by noise in small treatment groups, the hierarchical model smooths out some of that random variation.

5.2 *Who Responds to Alliances*

I used race, gender, hawkishness and internationalism to define demographic groups across which the impact of alliances might vary. I selected these variables because foreign policy dispositions like militant assertiveness shape general willingness to use force (Kertzer et al., 2014) as do gender (Barnhart et al., 2020) and race. I also control for other experimental manipulations.¹¹ Following TW's OLS analysis, I use a Gaussian likelihood, although the outcome is a binary variable.

I describe the results in two steps. First, I summarize the distribution of alliance effects. After this, I summarize the sources of variation in the alliance effect in Figure 3 and present the resulting heterogeneous effects for every group in Figure 4.

How alliances impact support for using force varies widely. Figure 2 provides an initial summary of that variation, and highlights several noteworthy estimates.

First, Figure 2 notes that the minimum estimated impact of an alliance on a demographic group is .05, while the maximum is .53. The maximum effect occurs among white men with high internationalism and low hawkishness. The minimum effect applies to non-white women with low internationalism and high hawkishness. There is no overlap in the posteriors of these estimates. The median group treatment effect estimate is .31, and this group of respondents is non-white men with middling internationalism and hawkishness. Alliances never clearly decrease support for intervention, but how much they increase support varies widely.

Figure 2 also presents the variation in how alliances impact demographic groups. The

¹¹See the appendix for priors.

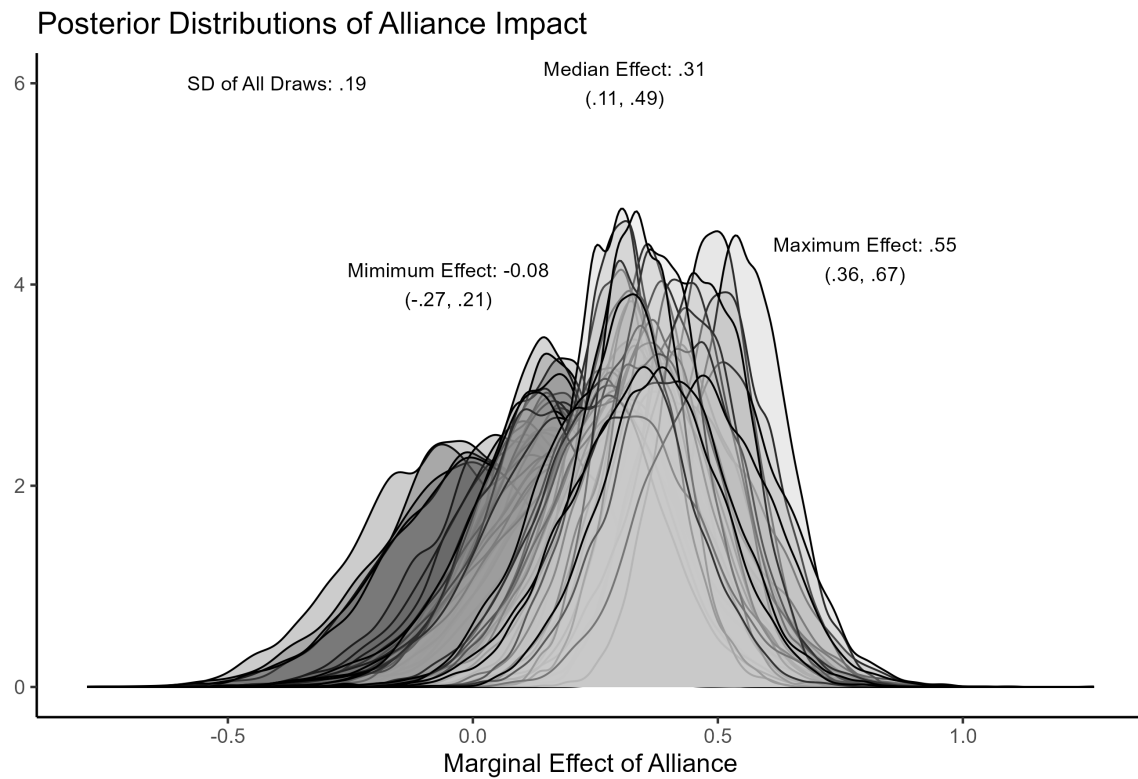


Figure 2. Posterior distribution of all estimated impacts of alliances on support for using force. Text values give notable point estimates, and parentheses summarize the 95% credible interval.

standard deviation of all posterior draws is .13. Roughly 5% of variation in the alliance effect is not explained by systematic regression components in Figure 3.¹²

Figure 3 plots how the impact of alliances varies across support for international engagement, willingness to use force, race and gender. These variables define the groups, so differences in the alliance slope follows their individual and joint variation. Because there are many groups, the impact of alliances varies widely within each level of these variables, but there are some clear patterns.

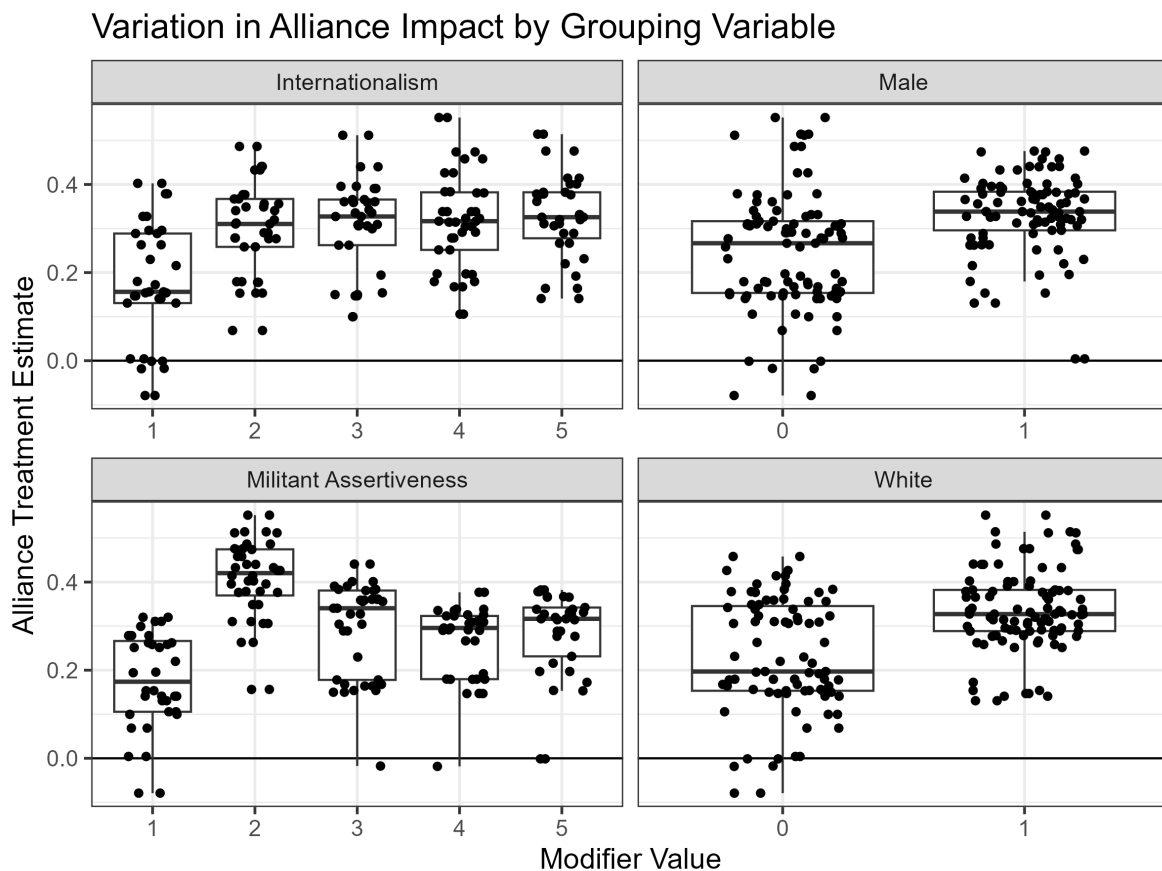


Figure 3. Variation in the impact of alliances on support for military intervention across four variables that set groups. Each point marks the impact of alliances on a specific group, and boxplots summarize the median and interquartile range of the slopes within each level of the variable. All slopes are present in each facet.

¹²This is σ_θ above.

Individuals with minimal interest in international engagement are less responsive to alliances, while any greater support for internationalism leads to a fairly consistent response to alliances. Similarly, alliances exert less impact on individuals who have minimal militant assertiveness. Alliances are very influential for individuals with low but greater than minimal hawkishness, however. Among those with moderate or high hawkishness, alliances have a fairly consistent impact. The media alliance impact is also greater for men, and among white respondents.

Crucially, parameteric interactions would not capture as clearly the non-linear steps across the levels of different grouping variables. Letting slopes vary across each level of the grouping variables generates more flexibility, and clearly shows the difference between individuals with very low internationalism or low militant assertiveness and others.

As Figure 3 suggests, alliances increase support for foreign intervention most among white men, especially those with low hawkishness and some internationalism. By contrast, alliances have little impact on support for war among non-white females who are also skeptical of international engagement and unwilling to use force. Individuals with more ambivalent foreign policy views respond more typically to TW's alliance treatment.

All these estimates suggest that internationalism matters more than hawkishness for understanding who is willing to fight for U.S. allies. Alliances may impact hawks less because these individuals support intervention regardless. Military alliances matter most to backers of international engagement who less willing to use force, but not entirely averse to military intervention.

These results show some of the strengths and weaknesses of the hierarchical approach to heterogeneous effects.¹³ A simple model based on demographic groups provides precise insights about who heeds alliances in supporting using force abroad. At the same time, because some demographic groups are small, the within-group effect estimates have substantial uncer-

¹³In the appendix, I analyze Bush and Prather (2020).

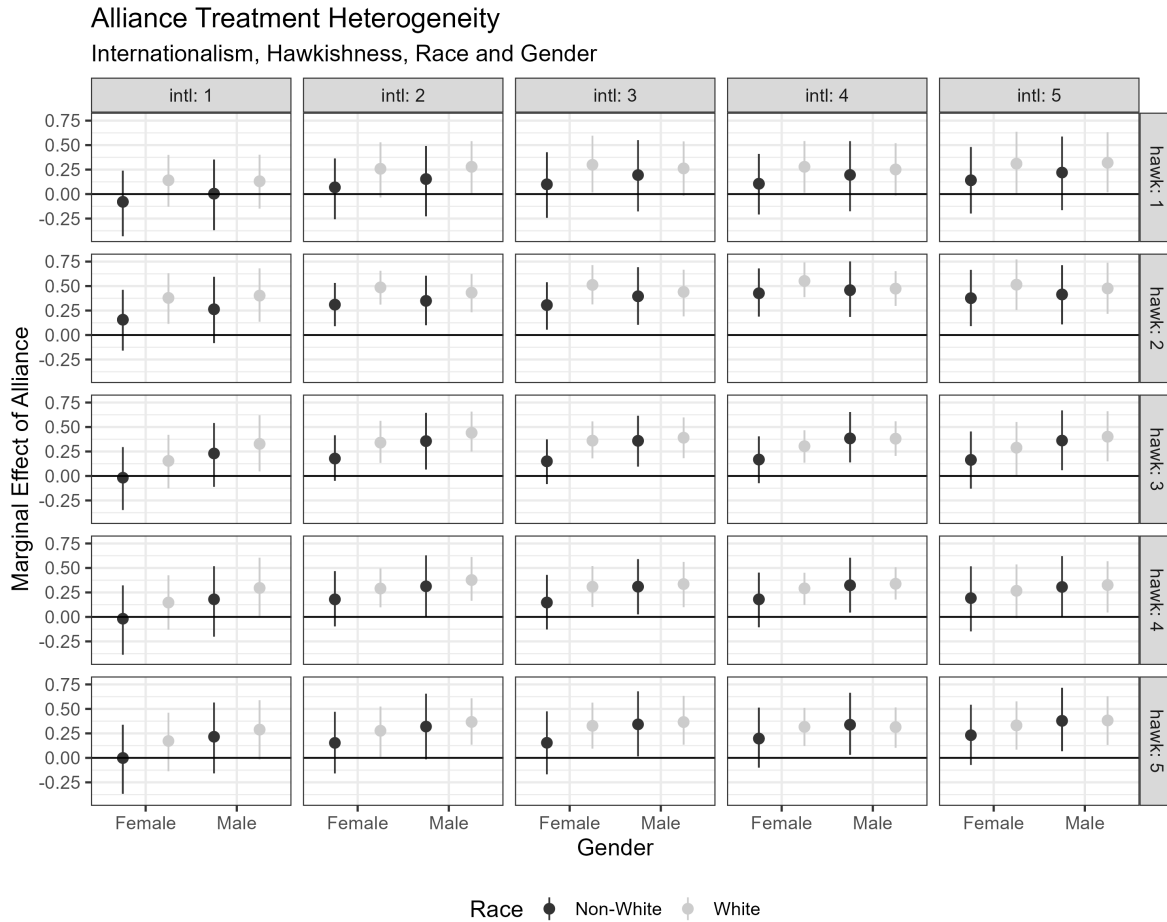


Figure 4. Estimates of the impact of military alliances on support for using force within demographic groups. Column facets are values of internationalism, and row facets are levels of hawkishness. X-axis divided by gender and colors demarcate gender. Points mark the posterior median and bars summarize the 95% credible interval.

tainty and powerful comparisons between most groups is challenging. Fewer groups would have more data and less uncertainty but perhaps obscure variation across key demographic characteristics.

6 Conclusion

This note explained how and when to use hierarchical models to estimate heterogeneous effects. Bayesian modeling can apply to a wide range of outcomes, data structures, and theories. It also details what drives variation in an effect and how much an effect varies. Explicitly modeling how different groups respond to an independent variable can help test arguments and inform policy.

Hierarchical modeling provides an intermediate approach between interactions or subgroup analyses and machine learning algorithms. For interactions with one or two variables, relying on simple interaction tools is best. Similarly, machine learning is best for discovery of complex heterogeneity. When there are two or more modifiers and many groups of theoretical interest, hierarchical modeling allows theoretically informed and interpretable estimation of effect variation.

As a result, hierarchical modeling complements existing tools and should not replace them. Researchers can use hierarchical models to check and inform other techniques, for instance by seeing if a key interaction holds when there are multiple modifiers, or comparing multiple modifiers that past theories have identified. Using hierarchical modeling can thus help scholars and policymakers better understand heterogeneous effects.

Acknowledgements

Thanks to Taylor Kinsley Chewning, Andrew Gelman and Carlisle Rainey for helpful comments.

References

- Abramson, Scott F, Korhan Koçak and Asya Magazinnik. 2022. “What Do We Learn about Voter Preferences from Conjoint Experiments?” *American Journal of Political Science* 66(4):1008–1020.
- Alley, Joshua. 2021. “Alliance Participation, Treaty Depth and Military Spending.” *International Studies Quarterly* 65(4):929–943.
- Arel-Bundock, Vincent. N.d. *marginaleffects: Predictions, Comparisons, Slopes, Marginal Means, and Hypothesis Tests*. R package version 0.14.0.9000.
URL: <https://vincentarelbundock.github.io/marginaleffects/>
- Arel-Bundock, Vincent, Ryan C Briggs, Hristos Doucouliagos, Marco Mendoza Aviña and Tom D Stanley. 2022. “Quantitative Political Science Research Is Greatly Underpowered.” *OSF Preprints*. Available at: <https://osf.io/preprints/osf/7vy2f>.
- Barnhart, Joslyn N, Robert F Trager, Elizabeth N Saunders and Allan Dafoe. 2020. “The Suffragist Peace.” *International Organization* 74(4):633–670.
- Blackwell, Matthew and Michael P Olson. 2022. “Reducing Model Misspecification and Bias in the Estimation of Interactions.” *Political Analysis* 30(4):495–514.
- Bürkner, Paul-Christian. 2017. “brms: An R package for Bayesian multilevel models using Stan.” *Journal of Statistical Software* 80(1):1–28.
- Bush, Sarah Sunn and Lauren Prather. 2020. “Foreign Meddling and Mass Attitudes Toward International Economic Engagement.” *International Organization* 74(2):584–609.
- Chaudoin, Stephen. 2014. “Promises or Policies? An Experimental Analysis of International Agreements and Audience Reactions.” *International Organization* 68(1):235–256.
- Clifford, Scott and Carlisle Rainey. 2023. Estimators for Topic-Sampling Designs. Technical report.
- Dorie, Vincent, George Perrett, Jennifer L Hill and Benjamin Goodrich. 2022. “Stan and BART for Causal Inference: Estimating Heterogeneous Treatment Effects Using the Power of Stan and the Flexibility of Machine Learning.” *Entropy* 24(12):1782.

- Feller, Avi and Andrew Gelman. 2015. "Hierarchical Models for Causal Effects." *Emerging Trends in the Social and Behavioral Sciences: An interdisciplinary, searchable, and linkable resource* pp. 1–16.
- Gelman, Andrew. 2008. "Scaling regression inputs by dividing by two standard deviations." *Statistics in medicine* 27(15):2865–2873.
- Gelman, Andrew. 2018. "You need 16 times the sample size to estimate an interaction than to estimate a main effect." Available at: <https://statmodeling.stat.columbia.edu/2018/03/15/need16/>.
- Goplerud, Max. 2021. "Modelling Heterogeneity Using Bayesian Structured Sparsity." *arXiv preprint arXiv:2103.15919*.
- Green, Donald P and Holger L Kern. 2012. "Modeling Heterogeneous Treatment Effects in Survey Experiments with Bayesian Additive Regression Trees." *Public Opinion Quarterly* 76(3):491–511.
- Grimmer, Justin, Solomon Messing and Sean J Westwood. 2017. "Estimating Heterogeneous Treatment Effects and the Effects of Heterogeneous Treatments with Ensemble Methods." *Political Analysis* 25(4):413–434.
- Guisinger, Alexandra and Elizabeth N. Saunders. 2017. "Mapping the Boundaries of Elite Cues: How Elites Shape Mass Opinion across International Issues." *International Studies Quarterly* 61(2):425–441.
- Imai, Kosuke and Marc Ratkovic. 2013. "Estimating treatment effect heterogeneity in randomized program evaluation." *The Annals of Applied Statistics* 7(1):443–470.
- Kertzer, Joshua D., Kathleen E. Powers, Brian C. Rathbun and Ravi Iyer. 2014. "Moral Support: How Moral Values Shape Foreign Policy Attitudes." *The Journal of Politics* 76(3):825–840.
- Kertzer, Joshua D and Ryan Brutger. 2016. "Decomposing Audience Costs: Bringing the Audience Back into Audience Cost Theory." *American Journal of Political Science* 60(1):234–249.
- Künzel, Sören R, Jasjeet S Sekhon, Peter J Bickel and Bin Yu. 2019. "Metalearners for estimating heterogeneous treatment effects using machine learning." *Proceedings of the national academy of sciences* 116(10):4156–4165.
- Levendusky, Matthew S and Michael C Horowitz. 2012. "When Backing Down is the Right Decision: Partisanship, New Information, and Audience Costs." *The Journal of Politics* 74(2):323–338.
- Marquardt, Kyle L. 2022. "Language, Ethnicity, and Separatism: Survey Results from Two Post-Soviet Regions." *British Journal of Political Science* 52(4):1831–1851.

- McElreath, Richard. 2016. *Statistical Rethinking: A Bayesian course with examples in R and Stan*. CRC Press.
- Schwartz, Joshua A and Christopher W Blair. 2020. “Do Women Make More Credible Threats? Gender Stereotypes, Audience Costs, and Crisis Bargaining.” *International Organization* 74(4):872–895.
- Simmons, Joseph P, Leif D Nelson and Uri Simonsohn. 2011. “False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant.” *Psychological Science* 22(11):1359–1366.
- Tomz, Michael and Jessica L.P. Weeks. 2021. “Military Alliances and Public Support for War.” *International Studies Quarterly* 65(3):811–824.
- Wager, Stefan and Susan Athey. 2018. “Estimation and Inference of Heterogeneous Treatment Effects using Random Forests.” *Journal of the American Statistical Association* 113(523):1228–1242.