# Distances in Latent Space: A Novel Approach to Analyzing Conjoints

Simon Hoellerbauer

July 14, 2020

# Motivation

- Often, we would like to place individuals in a latent space relative to other individuals or other fixed points within that space - politicians, tax plans, organizations, governments etc.

- We then want to know how the distance between an individual and such a fixed points affects that individual's attitudes and behaviors

- Often, the location of these points in latent space are influenced by a constellation of attributes - policy positions for politicians, tax rates and coverages for tax plans, structure, membership, goals for organizations, performance in different categories for governments

- This is really a **two-part process**, where closeness, decided by how individuals view an entity's traits, is the mechanism.

- I propose a methodological approach to studying this process

- Particularly well-suited to conjoint survey experiments, but could be adapted to a diverse array of experimental approaches.

# General Model

I describe a two part model:

1. IRT: to place individuals and profiles in the same latent space, making it possible to estimate the distance between them.
   - Based on random utility model, where individuals prefer profiles closer to them:

   $$U_{ij}(\mathbf{x}_j) = -(\theta_i - \xi(\mathbf{x}_j))^2 + \varepsilon_{ij}$$

   where $\xi(\mathbf{x}_j) = \mathbf{x}_j^\top \boldsymbol{\beta}$. $\mathbf{x}_j$ represents the vector of profile attributes for profile $j$.

2. Logistic Regression: to see how distance impacts a secondary outcome.

Note there are **separate outcome questions for each part**.

# IRT Portion of Model

$$
\begin{aligned}
\Pr(Y_{ik} = 1 | \mathbf{x}_{ik1}, \mathbf{x}_{ik2}) &= \Pr(U_{ik1} > U_{ik2}) \\
&= \Pr(-(\theta_i - \xi(\mathbf{x}_{ik1}))^2 + \varepsilon_{ik1} > -(\theta_i - \xi(\mathbf{x}_{ik2}))^2 + \varepsilon_{ik2}) \\
&= \Pr(-(\theta_i - \xi(\mathbf{x}_{ik1}))^2 + (\theta_i - \xi(\mathbf{x}_{ik2}))^2 > -\varepsilon_{ik1} + \varepsilon_{ik2}) \\
&= \Pr(-\theta_i^2 + \xi(\mathbf{x}_{ik1})\theta_i - \xi(\mathbf{x}_{ik1})^2 + \theta_i^2 - \xi(\mathbf{x}_{ik2})\theta_i + \xi(\mathbf{x}_{ik2})^2 > \varepsilon_{ik}) \\
&= \Pr(2(\xi(\mathbf{x}_{ik1}) - \xi(\mathbf{x}_{ik2}))\theta_i) + (-\xi(\mathbf{x}_{ik1})^2 + \xi(\mathbf{x}_{ik2})^2 > \varepsilon_{ik}) \\
&= \Pr(2(\xi(\mathbf{x}_{ik1}) - \xi(\mathbf{x}_{ik2}))\theta_i - (\xi(\mathbf{x}_{ik1})^2 + -\xi(\mathbf{x}_{ik2})^2) > \varepsilon_{ik}) \\
&= \Pr(2(\mathbf{x}_{ik1}^\top \boldsymbol{\beta} - \mathbf{x}_{ik2}^\top \boldsymbol{\beta})\theta_i - ((\mathbf{x}_{ik1}^\top \boldsymbol{\beta})^2 + -(\mathbf{x}_{ik2}^\top \boldsymbol{\beta})^2) > \varepsilon_{ik}) \\
&= \Phi(b(\mathbf{x}_{ik1}, \mathbf{x}_{ik2})\theta_i - g(\mathbf{x}_{ik1}, \mathbf{x}_{ik2}))
\end{aligned}
$$

If we assume $\varepsilon_{ik} \sim \mathcal{N}(0, \sigma)$, then $\Phi(.)$ represents the CDF of the Standard Normal distribution. This is then in the form of a two-parameter IRT model. $b(\mathbf{x}_{ik1}, \mathbf{x}_{ik2})$ and $g(\mathbf{x}_{ik1}, \mathbf{x}_{ik2})$ represent the item difficulty and combined item discrimination and item difficulty parameters, respectively where $b(\mathbf{x}_{ik1}, \mathbf{x}_{ik2}) = 2(x_{ik1} - x_{ik2})^\top \boldsymbol{\beta}/\sigma$ and $g(\mathbf{x}_{ik1}, \mathbf{x}_{ik2}) = \boldsymbol{\beta}^\top (\mathbf{x}_{ik1}\mathbf{x}_{ik1}^\top - \mathbf{x}_{ik2}\mathbf{x}_{ik2}^\top)\boldsymbol{\beta}/\sigma$.

# Logistic Regression Portion of Model

I connect the IRT model to the logistic regression via $\boldsymbol{\beta}$ and $\theta_i$, where I model the probability that an individual $i$ chooses profile 1 in profile pair $j$ or not (derived from the first outcome question listed above):

$$\Pr(W_{ij} = 1 | \mathbf{x}_{ij1}, \mathbf{x}_{ij2}) =$$
$$\text{logit}^{-1}(\gamma_0 + \gamma_1 * (2\theta_i(\mathbf{x}_{ij1} - \mathbf{x}_{ij2})^\top \boldsymbol{\beta} + \boldsymbol{\beta}^\top (\mathbf{x}_{ij2}\mathbf{x}_{ij2}^\top - \mathbf{x}_{ij1}\mathbf{x}_{ikj}^\top)\boldsymbol{\beta})$$

where $\boldsymbol{\beta}$ are the coefficients from the IRT model. Note that the term with the $\gamma_1$ coefficient is equal to $\theta_i - \mathbf{x}_{ij2}^\top\boldsymbol{\beta})^2 - (\theta_i - \mathbf{x}_{ij1}^\top\boldsymbol{\beta})^2$: : the *difference in the distance between ideal points and profile locations*: positive $= i$ closer to profile 1 than profile 2.
This is still a force-choice context; it is possible to adapt this approach in the case where a respondent faces separate choices for profile 1 and profile 2. Separate profile pairs used for each part of model.

Simon Hoellerbauer

# Application

- Conjoint Survey Experiment
- Research Question: How does the localness of organizations affect individual's willingness to interact with them?
- Project Goal 1: see if students feel closer to more local organizations - in the demographic sense and in the geographic sense.
- Project Goal 2: see if this closeness makes them more likely to declare a willingness to engage with an organization

## Application on Student Sample

- 676 students at University of North Carolina - Chapel Hill completed survey
- Each saw 15 profile-pairs, constructed from the following attribute-levels:

| Attribute | Level |
|---|---|
| Other members are | mainly students; students and non-students; mainly non-students |
| Leader is | a student; not a student |
| Organization's headquarters located in | Chapel Hill, NC; Raleigh, NC; Richmond, VA; Washington, DC |
| Organization is | not a chapter of a national organization; a chapter of a national organization |
| Funding mostly comes from | donations from members and community; donations from national partners |
| Aiming to increase voter registration | on campus; in the town of Chapel Hill |

# Application On Student Sample

- Respondents responded to 2 questions, always in the same order, after each pair:
    1. W: Would you be more likely to attend a meeting held by organization 1 or organization 2?
    2. Y: With which organization would you say you feel more of a personal connection?
- I used Y from profiles 2-15 for the IRT portion of the model
- I used W from profile 1 for the logistic regression portion of the model
- This was because of the possibility that the more profiles students saw, the more they would think about question 2 instead of question 1
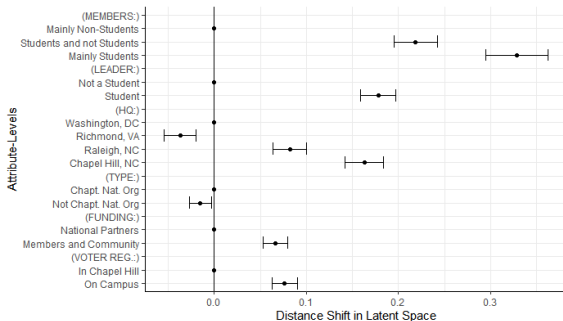
## Hypotheses:

1. Students will feel a greater affinity for student-involved and local organizations.
   1. Student-involved and local attribute-levels will place organizations to one side of the latent space.
   2. The mass of the ideal point distribution will be in the same portion of the latent space as all-student/all-local organizations.
2. An individual who is closer to organization 1 than organization 2 will be more likely to want to attend a meeting held by organization 1, and vice-version. In terms of the model, the coefficient on the difference in differences will be positive.

# Estimation

- Model fit using Stan
  - $\theta, \beta, \gamma \sim \mathcal{N}(0, 1)$
- For identification, $\theta$ was normalized to $\mathcal{N}(0, 1)$ and the coefficient on Leader: Student was fixed to be positive, to establish polarity of space.
- Traceplots and Rhat indicate that chains converged successfully

# Results: Student/Local Levels Consistently Place Organizations in Latent Space

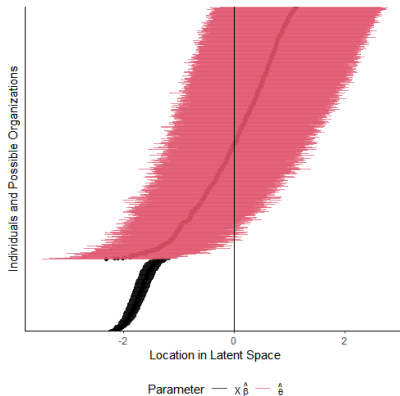Figure: $\hat{\beta}$—Determine Org. Locations (Posterior medians with 95% cred. int.)

# Results: Attribute-Level Coefficients Takeaways

- The most important attributes (largest coefficient size) represent demographically local organizational traits:
    - Identity of other members
    - Identity of leader
- Yet, geographic localness was also clearly important, with third largest coefficient on a Chapel Hill, NC headquarters
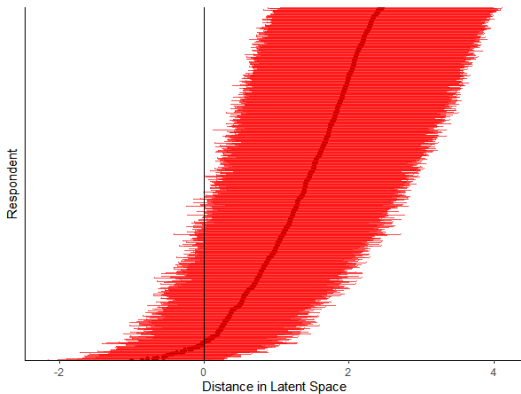- Local goals and local funding also mattered.

# Results: Most Respondents Closer to Student/Local Orgs

Figure: $\hat{\theta}$ and Possible Organization Positions ($\mathbf{X}\hat{\beta}$) (Posterior medians with 95% cred. int.)
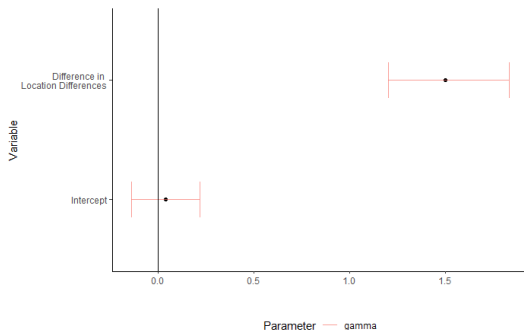
# Results: Most Respondents Closer to Student/Local Orgs

Figure: Difference Between Resp.'s Ideal Points and Most Student/Local Organization (Posterior medians with 95% cred. int.)
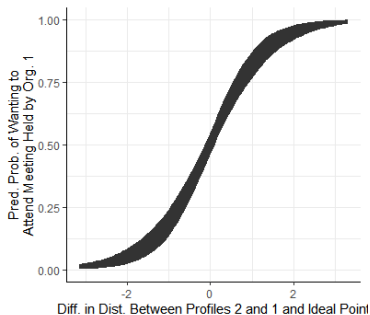
# Results: Resp. More Likely to Want to Attend Meeting of Org Closer to Them

Figure: Logistic Regression Model Coefficients Estimates ($\hat{\gamma}$) (Posterior medians with 95% cred. int.)

# Results: Resp. More Likely to Want to Attend Meeting of Org Closer to Them

Figure: Effect of Diff. in Dist. Between Ideal Points and Profiles on Probability of Wanting to Attend Meeting (Posterior medians with 95% cred. int.)

# Results: Summary

- Support found for both hypotheses
- Local organizational traits moved organizations closer to positive pole of latent space
- 59.3% of students had estimated positions with credible intervals entirely to the right of the most student/local organization
- As difference in distances increases—respondent is closer to org. 1 than org. 2—the probability of wanting to attend a meeting held by org. 1 increases (for reference, within data distances were normally distributed around 0, with standard deviation 1)

# Assessing Model Fit

- I used the part of the data that each portion of the model hadn't seen to assess out-of-sample prediction error.
- I use the area under the ROC curve (AUC). Because I have a sample of the posterior distribution of each parameter, I also can construct a picture of the AUC distribution.

Table: AUCs (Posterior Medians with 95% Credible Intervals)

| IRT: | 0.868 [0.858, 0.878] |
|---|---|
| Logistic Regression: | 0.765 [0.762, 0.769] |

- Logistic regression part of model does not fit as well; it is possible that students took other factors into account besides distance, or that the form of the distance is different (absolute difference, for example)

# Next Steps

- Restructure experiment so that respondents are asked only one question type after each profile pair
- Restructure experiment so that second question (W) is not forced-choice but asked about each profile in turn; this can get a better estimate of the effect of distance; requires modification of logistic regression portion of model
- Evaluate different distances in second part of model, not just squared distance
- Application for conjoints: perform typical conjoint AMCE analysis for W but also include distance
- More in-depth subgroup analysis
- Simulation study

# Appendix: Distribution of Difference in Distances

Figure: Differences in Distance Between Ideal Points and Profiles, Calculated Using $\hat{\theta}_i$ and $\hat{\beta}$